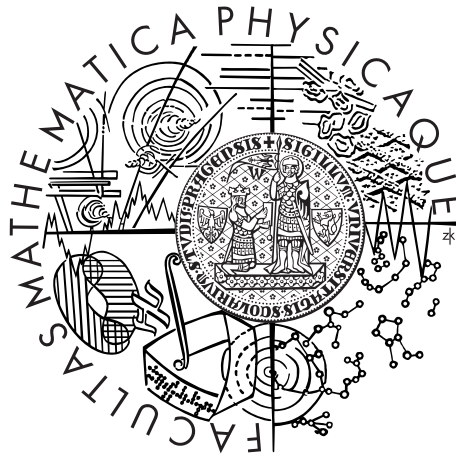


Univerzita Karlova v Praze  
Matematicko-fyzikální fakulta

## DIPLOMOVÁ PRÁCE



Lukáš Kotík

## Periodické regresní kvantily

Katedra (ústav): *Katedra pravděpodobnosti a matematické statistiky*  
Studijní program: *Matematika*  
Vedoucí diplomové práce: *RNDr. Daniel Hlubinka, Ph.D.*  
Studijní obor: *Pravděpodobnost, matematická statistika a ekonometrie*

Prohlašuji, že jsem svou diplomovou práci napsal samostatně a výhradně s použitím citovaných pramenů. Souhlasím se zapůjčováním práce.

V Praze dne

Lukáš Kotík

Název práce: *Periodické regresní kvantily*

Autor: *Lukáš Kotík*

Katedra (ústav): *Katedra pravděpodobnosti a matematické statistiky*

Vedoucí diplomové práce: *RNDr. Daniel Hlubinka, Ph.D.*

e-mail vedoucího: *hlubinka@karlin.mff.cuni.cz*

**Abstrakt:** *Práce se zabývá novým přístupem ke konstrukci konfidenčních množin pro vícerozměrné náhodné veličiny a pro vícerozměrné náhodné výběry. To lze také chápat jako jedno z možných rozšíření pojmu kvantil na více rozměrů.*

*Postup je založen na transformaci vycentrovaného náhodného vektoru do polárních (resp. hypersférických) souřadnic a poté určení tzv. směrových kvantilů. To jsou vlastně klasické jednorozměrné kvantily pro rozdělení poloměru podmíněné volbou úhlu polárních souřadnic. Výběrový protějšek směrového kvantilu odhadneme pomocí trigonometrické řady, jejíž koeficienty získáme kvantilovou regresí. Přejdem zpět ke kartézské soustavě souřadnic získáváme periodický regresní kvantil.*

*První kapitola je věnována volbě středu sloužícího k vycentrování dat. Nabídneme několik variant volby takového bodu. Zkoumána bude teoretická i výběrová varianta. Důraz bude především kladen na nejhlubší bod.*

*Druhá kapitola je věnována kvantilové regresi a zejména těm jejím rysům, které mají vliv na vlastnosti výběrového periodického regresního kvantilu.*

*Třetí a nejobsáhlejší kapitola již popisuje samotnou konstrukci a vlastnosti periodických regresních kvantilů. Popisována bude jak teoretická tak i výběrová varianta a jejich vzájemný vztah. Několik simulacních příkladů a zpracování reálných dat je uvedeno na konci této kapitoly.*

**Klíčová slova:** *konfidenční množiny, mnohorozměrný kvantil, nejhlubší bod, kvantilová regrese*

Title: *Periodical regression quantile*

Author: *Lukáš Kotík*

Department: *Department of Probability and Mathematical Statistics*

Supervisor: *RNDr. Daniel Hlubinka, Ph.D.*

Supervisor's e-mail address: *hlubinka@karlin.mff.cuni.cz*

**Abstract:** *The thesis deals with a new approach to construction of confidence regions for multivariate random variables and multivariate random samples. This can also be viewed as one of the possible generalizations of the notion of quantile into a multidimensional case.*

*The approach is based on the following: in the first step, a centred random vector is transformed into polar (hyperspherical) coordinates. Afterwards, so-called directional quantiles are determined. These are classical unidimensional quantiles for distribution of the radius conditional on the angle of the polar coordinates. Sample analogy of the directional quantiles is estimated using trigonometrical series with coefficients obtained by quantile regression.*

*The first chapter deals with the choice of the origin for the centralization of the data. We examine both theoretical and sample cases. We offer several variants with focus on the deepest point.*

*The second chapter concerns quantile regression with focus on the aspects, which have an impact on the properties of sample periodical regression quantiles.*

*The third and most exhaustive chapter is devoted to periodical regression quantiles construction and properties. Both theoretical and sample variants and their relationship are described. Several examples are offered in the end of the chapter.*

**Keywords:** *confidence regions, multidimensional quantile, deepest point, quantile regression*

# Obsah

<b>Úvod</b>	<b>4</b>
<b>1 Parametry polohy pro vícerozměrná data</b>	<b>5</b>
1.1 Nejhlubší bod . . . . .	5
1.2 Simplexový nejhlubší bod (Simplicial median) . . . . .	14
1.3 Mnohorozměrný median (Spatial median) . . . . .	16
1.4 Další možnosti zavedení mnohorozměrného mediánu . . . . .	16
<b>2 Kvantilová regrese</b>	<b>17</b>
<b>3 Periodické regresní kvantily</b>	<b>21</b>
3.1 Teoretické periodické regresní kvantily . . . . .	21
3.2 Výběrové periodické regresní kvantily . . . . .	27
Dvojměrný případ . . . . .	27
Vícerozměrný případ . . . . .	33
3.3 Metody pro zlepšení odhadu výběrových periodických kvantilů . . . . .	43
3.4 Volba řádu $p$ a konzistence výběrového periodického regresního kvantilu . . .	45
3.5 Pár poznámek . . . . .	46
3.6 Příklady výběrových kvantilů . . . . .	48
Exponenciální rozdělení – dvojměrný případ . . . . .	48
Normální rozdělení – dvojměrný případ . . . . .	52
Exponenciální rozdělení – trojměrný případ . . . . .	53
Normální rozdělení – trojměrný případ . . . . .	56
Příklad - porodní váha a délka . . . . .	57
<b>Literatura</b>	<b>59</b>

# Úvod

Jednorozměrné kvantily dnes bezesporu patří k základním charakteristikám popisu dat a náhodných veličin. Jejich uplatnění je velmi široké, namátkou zmiňme, že se např. používají při testech hypotéz, konstrukci konfidenčních intervalů, ale třeba také jako míra polohy jednorozměrných dat. K jejímu určení můžeme použít medián. Oproti často používanému průměru má tu výhodu, že odlehlá pozorování na něj nemají téměř žádný vliv. Ke konstrukci konfidenčního intervalu zase použijeme některý z krajních kvantilů. Zde je typické použití 97,5% a 2,5% kvantilů, které oddělují 2,5% největších a nejmenších hodnot.

V dnešní době se v mnoha situacích zaznamenává několik veličin vypovídajících o předmětu našeho zájmu. A tak by nás mohlo zajímat, zda je možné pojem kvantilu rozšířit na vícerozměrné náhodné veličiny a výběry. Pro ně se vlastně pojetí mnohorozměrného kvantilu a konfidenční množiny slévají.

Rozšíření kvantilu do více rozměrů není tak jednoduché a přináší sebou celou řadu otázek. Stejně jako jednorozměrný medián leží v jistém smyslu uprostřed rozložení dat na reálné ose, budeme i po mnohorozměrném mediánu vyžadovat, aby nějakým způsobem ležel uprostřed. Poznamenejme, že pojmy medián a 50% kvantil už pro vícerozměrná data nemusí být myšleno to samé. Zatímco první z těchto pojmů je chápán jako parametr polohy, druhý značí množinu, která „oddělí“ 50% krajních pozorování. A to je vlastně konfidenční množina.

Mediánu ve více rozměrech je dnes věnováno poměrně mnoho literatury. Většina definic je postavena na pojmu hloubka dat. Pomocí hloubky dat je také možné definovat mnohorozměrný kvantil. Jeho tvar však v mnoha případech není uspokojivý. Navíc je tento přístup velmi výpočetně náročný.

Tato práce přináší novou možnost zavedení mnohorozměrných kvantilů. Tento nový přístup je navíc poměrně přirozený a tedy i snadno interpretovatelný. Vychází z volby centrálního bodu a následného hledání klasických kvantilů ve všech směrech (přímkách) od tohoto bodu. Volbě těchto centrálních bodů je věnována první kapitola. Druhá kapitola se zabývá kvantilovou regresí jakožto nástroje, který použijeme k odhadu našeho kvantilu pro náhodné výběry. A konečně třetí kapitola popisuje samotnou konstrukci a vlastnosti jak teoretické, tak výběrové varianty *periodických regresních kvantilů*.

## Kapitola 1

# Parametry polohy pro vícerozměrná data

V této kapitole popíšeme několik možných zobecnění jednorozměrného mediánu pro vícerozměrné náhodné veličiny, jakožto parametru polohy. Stejně jako jednorozměrný medián dělí přímku na dvě části v kterých se bude vyskytovat přibližně polovina pozorování, je i pro tato vícerozměrná rozšíření žádoucí, aby v nějakém smyslu ležely uprostřed rozložení dat. Tyto body budou mít při konstrukci *periodických regresních kvantilů* velice důležitou úlohu. Budou totiž jakýmsi centrem konfidenčních množin vybudovaných naší metodou. Nabídneme více možností pro výběr tohoto centrálního bodu. Největší důraz však bude kladen na *nejhlubší bod*, který se pro tyto účely jeví jako nejvhodnější volba. Má mnoho pěkných vlastností, navíc je i dobře interpretovatelný.

### 1.1 Nejhlubší bod

Nejhlubší bod a poloprostorovou hloubku (halfspace depth) poprvé zavedl Tukey (1974). Halfspace depth má mnoho dobrých vlastností a je to zřejmě nejpoužívanější popis hloubky dat. Mezi její největší výhody patří její robustnost (viz Donoho & Gasko (1992)) a ekvivariance vzhledem k afinním transformacím. Nevýhodou je, že přímočarý výpočet nejhlubšího bodu je ve vyšších dimenzích a pro velké rozsahy výběru časově dosti náročný. Uspořádáme-li mnohorozměrná data podle jejich hloubky, dostaneme jakousi analogii jednorozměrného uspořádaného výběru. Bod s největší hloubkou můžeme považovat za zobecnění jednorozměrného mediánu.

**Definice 1.** *Mějme  $p$  rozměrný náhodný vektor  $\mathbf{X}$ . Nejhlubším bodem nazveme bod:*

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta} \in \mathbb{R}^p} \min_{\mathbf{u} \in \mathbb{R}^p} \mathbb{P}(\mathbf{u}^T \mathbf{X} \leq \mathbf{u}^T \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta} \in \mathbb{R}^p} \min_{\|\mathbf{u}\|=1} \mathbb{P}(\mathbf{u}^T \mathbf{X} \leq \mathbf{u}^T \boldsymbol{\theta})$$

*Číslo*

$$\text{depth}(\boldsymbol{\theta}) = \min_{\|\mathbf{u}\|=1} \mathbb{P}(\mathbf{u}^T \mathbf{X} \leq \mathbf{u}^T \boldsymbol{\theta})$$

*nazveme hloubkou bodu  $\boldsymbol{\theta}$ .*

**Definice 2.** Necht'  $\mathbf{X}_1, \dots, \mathbf{X}_n$  je náhodný výběr z nějakého  $p$  rozměrného rozdělení. Výběrovou hloubkou bodu  $\boldsymbol{\theta} \in \mathbb{R}^p$  nazveme hodnotu výrazu:

$$\text{depth}_n(\boldsymbol{\theta}) = \min_{\|\mathbf{u}\|=1} \text{card}\{i : \mathbf{u}^T \mathbf{X}_i \leq \mathbf{u}^T \boldsymbol{\theta}\}.$$

Výběrovým nejhlubším bodem (též Tukey median) nazveme bod:

$$\hat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \mathbb{R}^p} \text{depth}_n(\boldsymbol{\theta}).$$

Pro každý vektor  $\mathbf{u}$  a každý bod  $\boldsymbol{\theta}$  množina  $\{\mathbf{x} : \mathbf{u}^T \mathbf{x} \geq \mathbf{u}^T \boldsymbol{\theta}\}$  vymezuje poloprostor, který je oddělen nadrovinou procházející bodem  $\boldsymbol{\theta}$  a je kolmá na vektor  $\mathbf{u}$ . Ze všech nadrovin procházejících bodem  $\boldsymbol{\theta}$  vybereme tu, která vymezuje poloprostor s nejmenší počtem pozorování (nejmenší pravděpodobností výskytu pozorování). Pomocí ní pak určíme hloubku bodu. Nutno poznamenat, že jak by se mohlo zdát, není hloubka závislá jen na rozložení výskytu pozorování ve směru tohoto poloprostoru, ale na rozložení výskytu pozorování ve všech možných směrech od bodu  $\boldsymbol{\theta}$ . Je zřejmé, že čím více jsme na „okraji“ rozložení pozorování, tím je hloubka bodu menší. Naopak, vzdalujeme-li se od „okrajů“, hloubka bodů roste.

Označme  $H_{\mathbf{u}, \boldsymbol{\theta}}$  množinu  $\{\mathbf{y} : \mathbf{u}^T \mathbf{y} \leq \mathbf{u}^T \boldsymbol{\theta}\}$ . Pak zřejmě platí:

$$\hat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \mathbb{R}^p} \min_{\|\mathbf{u}\|=1} P_n(H_{\mathbf{u}, \boldsymbol{\theta}}) \quad (1.1)$$

$$\text{depth}_n(\boldsymbol{\theta}) = n \min_{\|\mathbf{u}\|=1} P_n(H_{\mathbf{u}, \boldsymbol{\theta}}), \quad (1.2)$$

kde  $P_n(A) = \frac{1}{n} \sum_{i=1}^n I\{\mathbf{X}_i \in A\}$  je empirická pravděpodobnostní míra.

Není-li nejhlubší bod určen jednoznačně (např. pro diskrétní rozdělení), budeme definovat nejhlubší bod jako průměr bodů s maximální hloubkou. Poznamenejme ještě, že i při jakémkoliv jiném „rozumném“ pravidlu pro určení jednoznačného nejhlubšího bodu (např. těžiště konvexního polyedru generovaného body s maximální hloubkou) zůstanou všechny důležité vlastnosti zachovány. Je zřejmé, že pro  $p = 1$  se bude nejhlubší bod shodovat s jednorozměrným mediánem.

Podívejme se nyní na nějaké užitečné vlastnosti nejhlubšího bodu.

**Věta 1.** Necht'  $p$  rozměrný náhodný vektor  $\mathbf{X}$  má spojitě rozdělení s hustotou  $f(\mathbf{x})$  a platí, že množina  $\mathcal{M} = \{\mathbf{x} : f(\mathbf{x}) > 0\}$  je souvislá. Pak existuje právě jeden nejhlubší bod.

*Důkaz.* Důkaz provedeme pro  $p = 2$ . Pro vyšší dimenze je postup analogický. Uvědomme si, že souvislost množiny  $\mathcal{M}$  implikuje následující vztah:

$$P(\mathbf{u}^T \mathbf{c} < \mathbf{u}^T \mathbf{X} < \mathbf{u}^T \mathbf{b}) > 0, \quad \forall \mathbf{u} \neq \mathbf{0}, \mathbf{c}, \mathbf{b} \in \mathcal{M} : \mathbf{u}^T \mathbf{c} < \mathbf{u}^T \mathbf{b} \quad (1.3)$$

Předpokládejme pro spor, že existují nejméně 2 nejhlubší body. Dva z nich si označme  $\mathbf{a}$  a  $\mathbf{b}$ . Dále si označme  $d = \text{depth}(\mathbf{a}) = \text{depth}(\mathbf{b})$  a  $\mathbf{t}$  vektor kolmý k přímce procházející body  $\mathbf{a}$ ,  $\mathbf{b}$ .

Uvažujme nejprve, že existuje bod  $\mathbf{o}$  ležící na úsečce  $\mathbf{ab}$  takový, že

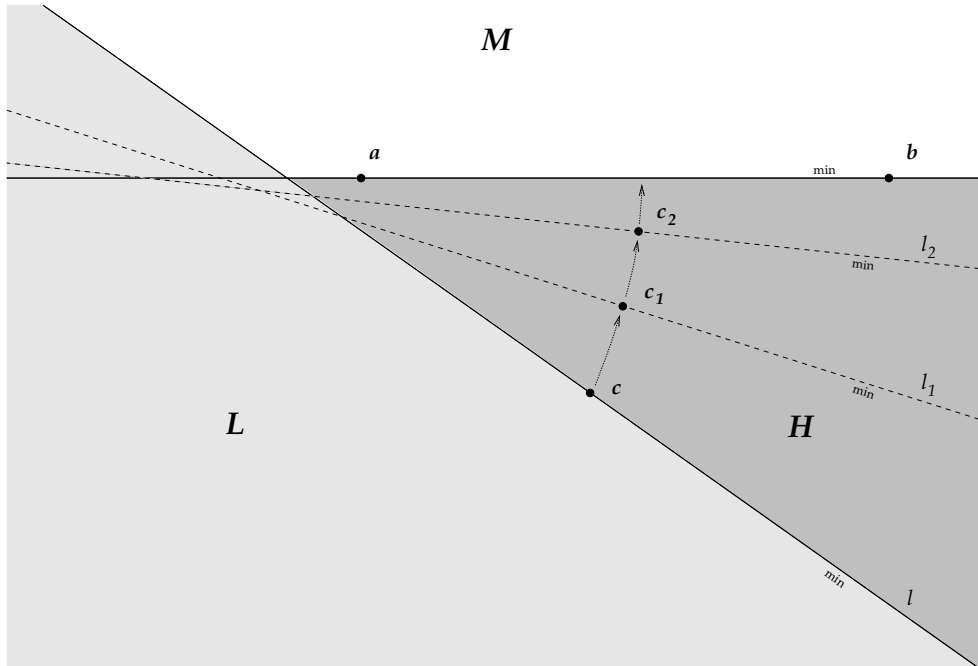
$$\mathbf{r} = \arg \min_{\|\mathbf{u}\|=1} P(\mathbf{u}^T \mathbf{X} \leq \mathbf{u}^T \mathbf{o})$$

není násobkem  $\mathbf{t}$  (tj. přímka vymežující polorovinu s nejmenší pravděpodobností procházející bodem  $\mathbf{c}$  není rovnoběžná s úsečkou  $\mathbf{ab}$ ). Platí  $\text{depth}(\mathbf{o}) \leq d$ . Zároveň však z (1.3) plyne, že buď  $P(\mathbf{r}^T \mathbf{X} \leq \mathbf{r}^T \mathbf{a}) < \text{depth}(\mathbf{o})$ , nebo  $P(\mathbf{r}^T \mathbf{X} \leq \mathbf{r}^T \mathbf{b}) < \text{depth}(\mathbf{o})$ . Tedy hloubka bodu  $\mathbf{a}$  nebo  $\mathbf{b}$  musí být ostře menší než  $d$ , tím se dostáváme ke sporu.

Pokud takový bod neexistuje, tj. pro libovolný bod  $\mathbf{o}$  z úsečky  $\mathbf{ab}$  existuje  $k \in \mathbb{R}$  tak, že platí

$$\mathbf{r} = \arg \min_{\|\mathbf{u}\|=1} P(\mathbf{u}^T \mathbf{X} \leq \mathbf{u}^T \mathbf{o}) = k\mathbf{t}.$$

Takže všechny body na úsečce  $\mathbf{ab}$  mají stejnou hloubku a stejnou polorovinu vymežující polo-prostor s minimální pravděpodobností procházející těmito body. Označme tento polo-prostor jako  $M$ . Pro libovolný bod  $\mathbf{c}$ , který neleží na přímce procházející body  $\mathbf{ab}$  a zároveň neleží v polorovině  $M$ , nastává situace vyobrazená na obr. 1.1: přímka  $l$  oddělující polorovinu s nejmenší pravděpodobností pro bod  $\mathbf{c}$  protíná přímku  $\mathbf{ab}$  vně úsečky  $\mathbf{ab}$  a tato polorovina neobsahuje úsečku  $\mathbf{ab}$ . V opačných případech bychom se podobně jako v předchozím odstavci, za pomoci vztahu (1.3), dostali ke sporu s hloubkou bodů  $\mathbf{a}$ ,  $\mathbf{b}$ . Tuto polorovinu si označme písmenem  $L$ . Dále, stejně jako na obr. 1.1, bude  $H$  značit oblast mezi přímkou  $l$  a přímkou procházející body  $\mathbf{a}$ ,  $\mathbf{b}$ .



Obrázek 1.1: Jednoznačnost nejhlubšího bodu. Symbol  $\text{min}$  označuje poloroviny s minimální pravděpodobností.

Vezměme posloupnost bodů  $\mathbf{c}_n$ ,  $n = 1, \dots, \infty$  ležících v  $H$  a vycházející z bodu  $\mathbf{c}$ , která bude konvergovat k nějakému bodu na úsečce  $\mathbf{ab}$ . Body  $\mathbf{c}_n$ ,  $n = 1, \dots, \infty$  budou mít stejné vlastnosti jako bod  $\mathbf{c}$ . Ke každému prvku posloupnosti zavedme, podobně jak tomu bylo pro bod  $\mathbf{c}$ , následující značení:  $l_i$  bude značit přímku procházející bodem  $\mathbf{c}_i$  a oddělující polorovinu  $L_i$  s nejmenší pravděpodobností pro tento bod. Úsek mezi  $l_i$  a přímkou procházející



body  $\mathbf{a}$  a  $\mathbf{b}$  označme  $H_i$ .

Z vlastností hloubky dostáváme:

$$P(L_n) \leq P(M) + P(H_n), \quad \forall n.$$

Z vlastností bodů mimo polorovinu  $M$  a díky faktu, že přímky mají nulovou míru, plyne:

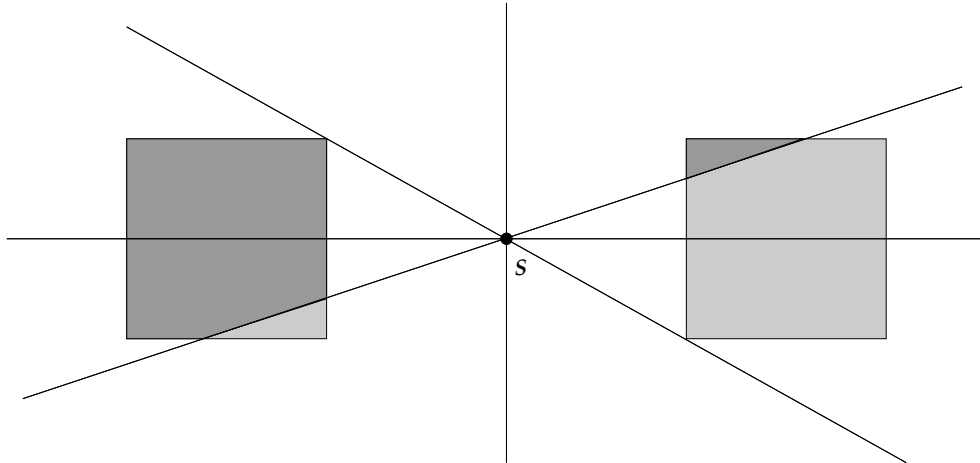
$$\lim_{n \rightarrow \infty} P(H_n) = 0, \quad \lim_{n \rightarrow \infty} P(L_n) = 1 - P(M).$$

Z toho dostáváme:

$$\lim_{n \rightarrow \infty} (P(L_n) - P(H_n)) = 1 - P(M) \leq P(M).$$

Takže  $P(M) = \frac{1}{2}$ . Všechny body na úsečce  $\mathbf{ab}$  mají tedy hloubku  $\frac{1}{2}$ . Libovolná přímka procházející těmito body, také musí dělit prostor na dvě poloroviny s pravděpodobnostní mírou  $\frac{1}{2}$ . Vezmeme-li jeden takový bod a libovolnou přímku různoběžnou s úsečkou  $\mathbf{ab}$ , dostaneme se díky (1.3) ke sporu s hloubkou ostatních bodů na úsečce  $\mathbf{ab}$ .  $\square$

**Poznámka.** Souvislost množiny  $M = \{\mathbf{x} : f(\mathbf{x}) > 0\}$  není nutnou podmínkou pro jednoznačnost nejhlubšího bodu. Představme si náhodný vektor s rovnoměrným rozdělením na dvou disjunktních čtvercích, tak jak je to zobrazeno na obr. 1.2. Pak bod  $S$  má hloubku rovnou  $\frac{1}{2}$ . Ostatní body mají hloubku ostře menší.



Obrázek 1.2: Jednoznačně určený nejhlubší bod pro rozdělení s nespojivou množinou  $\mathcal{M}$ .

**Věta 2.** Mějme  $p$  rozměrný náhodný vektor  $\mathbf{X}$  se spojitým rozdělením. Pak nejhlubší bod je ekvivalentní vůči libovolné afinní transformaci  $\mathbf{AX} + \mathbf{b}$ , kde  $\mathbf{A} \in \mathbb{R}^{p \times p}$  je regulární matice a  $\mathbf{b} \in \mathbb{R}^p$ .

Jinými slovy nejhlubší bod náhodného vektoru  $\mathbf{AX} + \mathbf{b}$  je roven  $\mathbf{A}\hat{\boldsymbol{\theta}} + \mathbf{b}$ , kde  $\hat{\boldsymbol{\theta}}$  značí nejhlubší bod vektoru  $\mathbf{X}$ .

Důkaz.

$$\max_{\boldsymbol{\xi} \in \mathbb{R}^p} \min_{\mathbf{v} \in \mathbb{R}^p} P(\mathbf{v}^T \mathbf{AX} \leq \mathbf{v}^T \boldsymbol{\xi}) = \max_{\boldsymbol{\xi} \in \mathbb{R}^p} \min_{\mathbf{v} \in \mathbb{R}^p} P(\mathbf{v}^T \mathbf{AX} \leq \mathbf{v}^T \mathbf{A} \mathbf{A}^{-1} \boldsymbol{\xi}) \quad (1.4)$$

$$= \max_{\boldsymbol{\theta} \in \mathbb{R}^p} \min_{\mathbf{u} \in \mathbb{R}^p} P(\mathbf{u}^T \mathbf{X} \leq \mathbf{u}^T \boldsymbol{\theta}), \quad (1.5)$$

kde  $\mathbf{u} = \mathbf{v}^T \mathbf{A}$  a  $\boldsymbol{\theta} = \mathbf{A}^{-1} \boldsymbol{\xi}$ . Regulárnost matice  $\mathbf{A}$  nám zaručuje, že afinní zobrazení dané maticemi  $\mathbf{A}$  a  $\mathbf{A}^{-1}$  je prosté a na  $\mathbb{R}^p$ . Proto je přechod od (1.4) k (1.5) oprávněný. Výraz v (1.5) nabývá minima pro  $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}} = \mathbf{A}^{-1} \hat{\boldsymbol{\xi}}$ , kde  $\hat{\boldsymbol{\xi}}$  je výraz, který minimalizuje výraz na levé straně v (1.4).  $\hat{\boldsymbol{\xi}} = \mathbf{A} \hat{\boldsymbol{\theta}}$  je tedy nejhlubším bodem rozdělení náhodného vektoru  $\mathbf{A}\mathbf{X}$ . Ekvivariance vůči posunutí o libovolný vektor  $\mathbf{b}$  je zřejmá.  $\square$

Z důkazu věty 2 je vidět, že také hloubka je invariantní vůči afinním transformacím.

**Věta 3.** *Výběrový nejhlubší bod (Tukey median) je ekvivantní vůči libovolné afinní transformaci  $\mathbf{A}\mathbf{X} + \mathbf{b}$ , kde  $\mathbf{A} \in \mathbb{R}^{p \times p}$  je regulární matice a  $\mathbf{b} \in \mathbb{R}^p$ .*

*Důkaz.* V důkazu věty 2 stačí místo pravděpodobnostní míry  $P$  použít  $P_n$ .  $\square$

**Věta 4.** *Nechť  $\mathbf{X}_1, \dots, \mathbf{X}_n$  jsou nezávislé,  $\mathbf{X}_i \sim P$ ,  $i = 1, \dots, n$ . Pak pro libovolný bod  $\mathbf{x}$  platí:*

$$\frac{1}{n} \text{depth}_n(\mathbf{x}) \xrightarrow{n \rightarrow \infty} \text{depth}(\mathbf{x}) \quad \text{s.j.}$$

*Důkaz.* Stačí si uvědomit, že platí (použijeme definici výběrové hloubky pomocí vztahu (1.2)):

$$\sup_{\boldsymbol{\theta}} |n^{-1} \text{depth}_n(\boldsymbol{\theta}) - \text{depth}(\boldsymbol{\theta})| \leq \sup_{\mathbf{u}, \boldsymbol{\theta}} |P_n(H_{\mathbf{u}, \boldsymbol{\theta}}) - P(H_{\mathbf{u}, \boldsymbol{\theta}})|$$

a pro  $\mathbf{X}_i$  *i.i.d.*  $P$  platí:

$$\sup_A |P_n(A) - P(A)| \xrightarrow{n \rightarrow \infty} 0 \quad P \text{ s.j.}$$

$\square$

Z věty 4 okamžitě plyne:

$$\hat{\boldsymbol{\theta}}_n \xrightarrow{n \rightarrow \infty} \hat{\boldsymbol{\theta}} \quad \text{s.j.}$$

Jak ukáže následující věta, nejhlubší bod je poměrně robustním odhadem.

**Věta 5.** *Bod zlomu výběrového nejhlubšího bodu je nejméně roven  $\frac{1}{p+1}$ .*

*Důkaz.* Důkaz je uveden v Donoho & Gasko (1992).  $\square$

Ve skutečnosti jsme většinou od této hranice poměrně vzdáleni. Dá se ukázat, že pro náhodný výběr ze středově symetrického rozdělení, tj. platí  $P(\mathbf{a} + S) = P(\mathbf{a} - S)$ , pro všechny měřitelné množiny  $S$  a nějaký bod  $\mathbf{a}$  (střed symetrie), konverguje bod zlomu nejhlubšího bodu pro  $p > 2$  s rostoucím  $n$  s.j. k  $1/3$ .

Jakou hloubku můžeme očekávat u nejhlubšího bodu? V Donoho & Gasko (1992) je ukázáno, že maximální hloubka leží mezi  $\lceil n/(p+1) \rceil$  a  $\lceil n/2 \rceil$ . Přičemž pro středově symetrická rozdělení konverguje  $n^{-1} \text{depth}_n(\hat{\boldsymbol{\theta}}_n)$  s.j. k  $1/2$ .

### Příklad 1 Středově symetrické rozdělení

Uvažujme rozdělení, které splňují předpoklady věty 1. Sem např. patří i mnohorozměrné normální rozdělení. Nejhlubším bodem takto rozdělených veličin je střed symetrie. BÚNO uvažujme středovou symetrii kolem  $\mathbf{0}$ . Pokud by nejhlubší bod, označme si ho  $\mathbf{a}$ , byl různý od středu symetrie, pak by ale i bod  $-\mathbf{a}$  byl také nejhlubší a to je spor s tvrzením věty 1. Ze středové symetrie také dostáváme, že každá nadrovina procházející středem, dělí prostor

na dva podprostory se stejnou pravděpodobností (rovné 1/2). Nejhlubší bod má tedy v tomto případě hloubku 1/2.

Pokud by nejhlubší bod nebyl určen jednoznačně (např. středově symetrické diskrétní rozdělení), pak ke každému bodu s maximální hloubkou existuje středově symetrický bod se stejnou hloubkou. Předefinujeme-li nejhlubší bod jako těžiště konvexního polyedru generovaného těmito body, pak se zřejmě bude shodovat se středem symetrie.

### Příklad 2 Exponenciální rozdělení

Mějme náhodný vektor  $\mathbf{X} = (X_1, X_2)$ .  $X_1, X_2$  nezávislé,  $X_i \sim \text{Exp}(\lambda_i)$  ( $f_i(x) = \lambda_i e^{-\lambda_i x}$ ),  $i = 1, 2$ . Zkusme spočítat nejhlubší bod.

Nejprve si uvědomme, že pokud  $Z \sim \text{Exp}(1)$ , potom náhodná veličina  $Y = \frac{1}{\lambda}Z$  má exponenciální rozdělení s parametrem  $\lambda$ . Stačí tedy spočítat nejhlubší bod pro náhodný vektor jehož složky mají rozdělení s parametrem 1. Nejhlubší bod pro vektor  $\mathbf{X}$  dostaneme tudíž jako

$$\hat{\xi} = \begin{pmatrix} 1/\lambda_1 & 0 \\ 0 & 1/\lambda_2 \end{pmatrix} \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix},$$

kde  $(\hat{\theta}_1, \hat{\theta}_2)^T$  je nejhlubší bod rozdělení náhodného vektoru se složkami s exponenciálním rozdělením s parametry rovnými jedné. Ještě uveďme jedno pomocné tvrzení.

**Lemma.** *Teoretický nejhlubší bod dvojrozměrného náhodného vektoru  $\mathbf{Z}$  s nezávislými složkami s exponenciálním rozdělením s parametrem 1 leží na přímce  $z_2 = z_1$ .*

*Důkaz.* Předpokládejme, že nejhlubší bod  $(\theta_1, \theta_2)^T$  leží mimo tuto přímku. Ze symetrie rozdělení vektoru  $\mathbf{Z}$  kolem osy  $z_2 = z_1$  dostáváme, že bod  $(\theta_2, \theta_1)^T$  musí mít stejnou hloubku jako  $(\theta_1, \theta_2)^T$ . Je tedy také nejhlubší. To je spor s tvrzením věty 1.  $\square$

Stačí tedy počítat hloubku jen pro body ležící na přímce  $z_2 = z_1$ . Označme si  $\theta = \theta_1 = \theta_2$ . Použijeme jinou parametrizaci: místo  $P(\mathbf{u}^T \mathbf{Z} \leq \mathbf{u}^T \boldsymbol{\theta})$  budeme počítat  $P(Z_2 \leq aZ_1 + b)$ , kde přímka  $z_2 = az_1 + b$  prochází bodem  $(\theta, \theta)^T$ . Takže platí  $\theta = a\theta + b$  a tedy  $b = \theta(1 - a)$ . Zajímá nás tedy výpočet  $P(Z_2 \leq aZ_1 + \theta(1 - a))$ . Dále platí:

$$\min_{\mathbf{u} \in \mathbb{R}^2} P(\mathbf{u}^T \mathbf{Z} \leq \mathbf{u}^T \boldsymbol{\theta}) = \min_{a \in \mathbb{R}} \min\{P(Z_2 \leq aZ_1 + \theta(1 - a)), P(Z_2 \geq aZ_1 + \theta(1 - a))\} \quad (1.6)$$

Výpočet  $P(Z_2 \leq aZ_1 + \theta(1 - a))$  provedeme zvlášť pro různé hodnoty parametru  $a$ :

1.  $0 \leq a \leq 1$

$$\begin{aligned}
P(Z_2 \leq aZ_1 + \theta(1-a)) &= \int_0^\infty \left( \int_0^{az_1 + \theta(1-a)} e^{-z_2} dz_2 \right) e^{-z_1} dz_1 \\
&= \int_0^\infty \left( 1 - e^{-(az_1 + \theta(1-a))} \right) e^{-z_1} dz_1 \\
&= 1 - e^{-\theta(1-a)} \int_0^\infty e^{-z_1(a+1)} dz_1 \\
&= 1 - \frac{e^{-\theta(1-a)}}{1+a}.
\end{aligned}$$

2.  $1 < a$

$$\begin{aligned}
P(Z_2 \leq aZ_1 + \theta(1-a)) &= P(Z_2 \geq \frac{1}{a}Z_1 - \frac{\theta(1-a)}{a}) \\
&= \frac{a \exp\{\theta(1-a)/a\}}{1+a}.
\end{aligned}$$

V poslední rovnosti jsme využili výsledku z předchozí části.

3.  $a < 0$

$$\begin{aligned}
P(Z_2 \leq aZ_1 + \theta(1-a)) &= \\
&= \int_0^{-\theta(1-a)/a} \left( \int_0^{az_1 + \theta(1-a)} e^{-z_2} dz_2 \right) e^{-z_1} dz_1 \\
&= \int_0^{-\theta(1-a)/a} e^{z_1} dz_1 - e^{-\theta(1-a)} \int_0^{-\theta(1-a)/a} e^{-z_1(a+1)} dz_1 \\
&= \begin{cases} 1 - \frac{1}{1+a} (a e^{\theta(1-a)/a} + e^{-\theta(1-a)}), & a \neq -1 \\ 1 - e^{-2\theta}(1+2\theta), & a = -1 \end{cases}
\end{aligned}$$

Označme  $p(a, \theta) = P(Z_2 \geq aZ_1 + \theta(1-a))$ . Pak:

$$p(a, \theta) = \begin{cases} \frac{1}{1+a} (a e^{\theta(1-a)/a} + e^{-\theta(1-a)}), & a < 0, a \neq -1 \\ e^{-2\theta}(1+2\theta), & a = -1 \\ \frac{e^{-\theta(1-a)}}{1+a}, & 0 \leq a \leq 1 \\ 1 - \frac{a \exp\{\theta(1-a)/a\}}{1+a}, & 1 < a \end{cases}$$

Jak je vidět z (1.6), bude nás zajímat průběh následující funkce:

$$g(a, \theta) = \min \{p(a, \theta), 1 - p(a, \theta)\},$$

pro kterou platí:

$$\text{depth}_P((\theta, \theta)^T) = \min_{a \in \mathbb{R}} g(a, \theta).$$

Pro nalezení nejhlubšího bodu musíme nejdříve najít minimum funkce  $g(a, \theta)$  v parametru  $a$ . Pro větší přehlednost se nejdříve budeme zabývat plochou nad přímkou, tj. funkcí  $p(a, \theta)$ . Snadno se zjistí, že je to funkce spojitá v proměnné  $a$  pro libovolné  $\theta > 0$ .

$$\frac{\partial p(a, \theta)}{\partial a} = \begin{cases} \left( (-\theta a^{-1} - \theta \frac{1-a}{a^2}) e^{\theta \frac{1-a}{a}} a + e^{\theta \frac{1-a}{a}} + \theta e^{-\theta+\theta a} \right) (a+1)^{-1} \\ \quad - \left( e^{\theta \frac{1-a}{a}} a + e^{-\theta+\theta a} \right) (a+1)^{-2}, & a < 0, a \neq -1 \\ 0, & a = -1 \\ \theta e^{-\theta+\theta a} (a+1)^{-1} - e^{-\theta+\theta a} (a+1)^{-2}, & 0 \leq a \leq 1 \\ -e^{-\theta+\theta a^{-1}} (a+1)^{-1} + \theta e^{-\theta+\theta a^{-1}} a^{-1} (a+1)^{-1} \\ \quad + a e^{-\theta+\theta a^{-1}} (a+1)^{-2}, & a > 1 \end{cases}$$

Derivace je spojitá funkce a je rovna nule pro  $a \in \left\{-1, \frac{1-\theta}{\theta}, \frac{\theta}{1-\theta}\right\}$ . Podívejme se nyní podrobněji na průběh funkce  $p(a, \theta)$  v okolí těchto bodů. Nejhlubší bod by měl ležet někde mezi 0.5 až 1, omezíme se proto při vyšetřování průběhu funkce jen na tyto hodnoty parametru  $\theta$ :

1.  $a = -1$

Dalšími výpočty zjistíme, že

$$\lim_{a \rightarrow -1} \frac{\partial^2 p(a, \theta)}{\partial a^2} = \frac{2\theta^3 - 3\theta^2}{3e^{2\theta}}$$

Tento výraz je záporný pro  $\theta \in (0, 1.5)$ . Takže bod -1 je bodem lokálního maxima a platí  $p(-1, \theta) = e^{-2\theta}(1 + 2\theta)$

2.  $a = \frac{1-\theta}{\theta}$

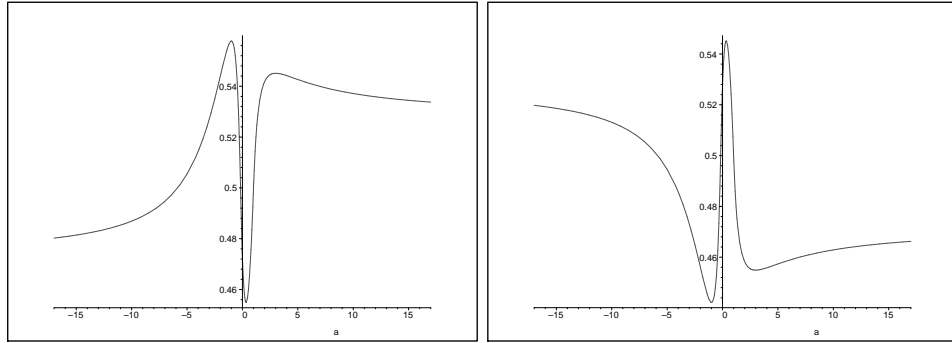
Pro  $\theta \in (0.5, 1)$  leží bod  $\frac{1-\theta}{\theta}$  v intervalu  $(0, 1)$ . Pro  $0 < a < \frac{1-\theta}{\theta}$  je derivace záporná, pro  $a$  za tímto bodem je kladná. Máme tedy bod lokálního minima a dostáváme  $p\left(\frac{1-\theta}{\theta}, \theta\right) = \theta e^{1-2\theta}$ .

3.  $a = \frac{\theta}{1-\theta}$

Podobně jako v předchozím bodu zjistíme, že  $\frac{\theta}{1-\theta}$  je bodem lokálního maxima a  $p\left(\frac{\theta}{1-\theta}, \theta\right) = 1 - \theta e^{1-2\theta}$

Podívejme se ještě na limity v krajních bodech:

$$\begin{aligned} \lim_{a \rightarrow \infty} p(a, \theta) &= 1 - e^{-\theta} \\ \lim_{a \rightarrow -\infty} p(a, \theta) &= e^{-\theta} \end{aligned}$$

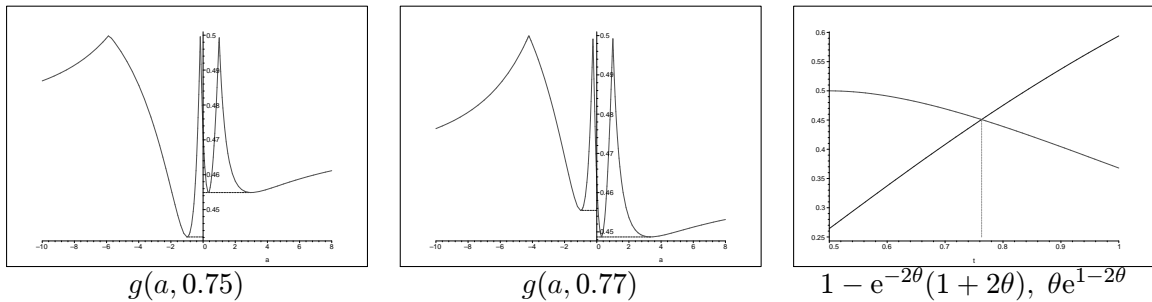
Obrázek 1.3: Obsah plochy pod a nad přímkou  $az_1 + \theta(1 - a)$ 

Vlastnosti funkce  $1 - p(a, \theta)$  získáme snadno z vlastností funkce  $p(a, \theta)$ . ( $1 - p(a, \theta)$  je s  $p(a, \theta)$  osově souměrná podél přímky  $y = 1/2$ ). Na obr. 1.3 je k vidění průběh  $p(a, \theta)$  a  $1 - p(a, \theta)$  jako funkcí proměnné  $a$ .

Nás však zajímá funkce  $g(a, \theta)$  a také hodnoty proměnné  $a$  v nichž nabývá svého minima. Z výše napsaného plyne, že to budou body:  $-1, \frac{1-\theta}{\theta}, \frac{\theta}{1-\theta}, -\infty, \infty$ .

Platí  $\lim_{a \rightarrow -\infty} g(a, \theta) = \lim_{a \rightarrow +\infty} g(a, \theta) = e^{-\theta} = g(0, \theta)$  a vzhledem k tomu, že funkce  $g(\cdot, \theta)$  na intervalu  $(0, \frac{1-\theta}{\theta})$  klesá tak se o body  $-\infty, +\infty$  nemusíme zajímat.

Dále  $g(\frac{1-\theta}{\theta}, \theta) = g(\frac{\theta}{1-\theta}, \theta) = \theta e^{1-2\theta}$  a  $g(-1, \theta) = 1 - e^{-2\theta}(1 + 2\theta)$ . Chceme najít bod globálního minima funkce  $g(\cdot, \theta)$ . Z obrázku 1.4 je vidět, že zřejmě pro hodnoty  $\theta$  menší než nějaké  $\theta_0$  bude  $g$  nabývat minima pro  $a = -1$  a pro ostatní hodnoty  $\theta$  v  $a = \frac{1-\theta}{\theta}, \frac{\theta}{1-\theta}$ .

Obrázek 1.4:  $g(a, \theta)$  pro dvě různé hodnoty  $\theta$  a velikosti lokálních (globálních) extrémů funkce  $g$ 

Velikosti lokálních minim  $g$  jsou rovny  $1 - e^{-2\theta}(1 + 2\theta)$  a  $\theta e^{1-2\theta}$ . První z nich je funkcí rostoucí (připomeňme, že se zajímáme jen o interval  $(0.5, 1)$ ), druhá funkcí klesající. Je zřejmé, že globální minimum funkce  $g(a, \theta)$  v proměnné  $a$  je rovno  $\min \{1 - e^{-2\theta}(1 + 2\theta), \theta e^{1-2\theta}\}$ . Označme  $\theta_0$  bod v němž se obě funkce protínají, pak  $\min \{1 - e^{-2\theta}(1 + 2\theta), \theta e^{1-2\theta}\}$  je

pro  $\theta \leq \theta_0$  rostoucí a pro  $\theta \geq \theta_0$  klesající. Pro bod  $\theta_0$  tedy platí:

$$\begin{aligned}\theta_0 &= \arg \max_{\theta} \min \left\{ 1 - e^{-2\theta}(1 + 2\theta), \theta e^{1-2\theta} \right\} \\ &= \arg \max_{\theta} \min_{a \in \mathbb{R}} \min \{ P(Z_2 \leq aZ_1 + \theta(1 - a)), P(Z_2 > aZ_1 + \theta(1 - a)) \} \\ &= \arg \max_{\theta=(\theta, \theta)^T} \min_{\mathbf{u} \in \mathbb{R}^2} P(\mathbf{u}^T \mathbf{Z} \leq \mathbf{u}^T \boldsymbol{\theta})\end{aligned}$$

Teoretickým nejhlubším bodem je tedy bod  $(\theta_0, \theta_0)^T$ , kde  $\theta_0$  je řešením rovnice:

$$1 - e^{-2\theta}(1 + 2\theta) = \theta e^{1-2\theta} \quad (1.7)$$

Rovnici (1.7) nelze přesně spočítat. Numerickým řešením rovnice dostaneme

$$\begin{aligned}\theta_0 &\doteq 0.7630687272, \\ \text{depth}((\theta_0, \theta_0)^T) &\doteq 0.45088.\end{aligned}$$

Vidíme, že ani výpočet teoretického nejhlubšího bodu není triviální záležitostí.

Přímočarý výpočet výběrového nejhlubšího bodu je velice časově náročný. V 90. letech minulého století se objevilo několik přesných i aproximativních algoritmů, založených na geometrických vlastnostech hloubky, pro rychlejší výpočet nejhlubšího bodu.

Pro výpočet přesné hodnoty nejhlubšího bodu můžeme využít algoritmy HALFMED, LDEPTH a ISODEPTH. HALFMED a ISODEPTH fungují jen pro dvojrozměrná data, LDEPTH pro  $p > 3$  najde jen aproximaci nejhlubšího bodu. HALFMED spočítá výběrový nejhlubší bod v čase  $O(n^2 \log^2 n)$ . LDEPTH a ISODEPTH jsou pomalejší, ale pořád dosahují výrazně lepších výsledků než přímočarý algoritmus ( $O(n^5 \log n)$ ).

V Struyf & Rousseeuw (2000) je popsán algoritmus DEEPLOC, který dovede spočítat aproximaci nejhlubšího bodu pro  $p$  rozměrná data v čase  $O(5n^{0.3}(mpn \log n + p^2n) + mp^3 + mpn)$ , kde  $m$  je konstanta, která značí počet směrů pomocí kterých algoritmus počítá nejhlubší bod. Většinou se volí  $m = 500$ .

Pro více než trojrozměrná data můžeme tedy pro výpočet aproximace nejhlubšího bodu použít algoritmy DEEPLOC a LDEPTH. DEEPLOC je však podstatně rychlejší. U obou těchto algoritmů platí, že s rostoucím  $n$  konverguje jimi spočtená aproximace ke skutečné hodnotě výběrového nejhlubšího bodu.

Na <http://www.agoras.ua.ac.be/Locdept.htm> je k stáhnutí implementace těchto algoritmů do fortranu a jejich podrobný popis.

Na <http://www.r-project.org/> se dá stáhnout funkce *bagplot.R*, která mimo jiné umí spočítat nejhlubší bod (pomocí algoritmu HALFMED) pro dvojrozměrná data.

## 1.2 Simplexový nejhlubší bod (Simplicial median)

Další možnost zavedení hloubky dat se dá nalézt v Liu (1990). Podívejme se na definici a nějaké základní vlastnosti simplexové hloubky.

**Definice 3.** Mějme  $\mathbf{X}_1, \dots, \mathbf{X}_{p+1}$  i.i.d. na  $\mathbb{R}^p$ . Označme  $S[\mathbf{X}_1, \dots, \mathbf{X}_{p+1}]$  simplex s vrcholy  $\mathbf{X}_1, \dots, \mathbf{X}_{p+1}$  (tj.  $S[\mathbf{X}_1, \dots, \mathbf{X}_{p+1}]$  je množina všech konvexních kombinací  $\mathbf{X}_1, \dots, \mathbf{X}_{p+1}$ ). Pak simplexovou hloubkou bodu  $\mathbf{x} \in \mathbb{R}^p$  rozumíme:

$$SD(\mathbf{x}) = P(\mathbf{x} \in S[\mathbf{X}_1, \dots, \mathbf{X}_{p+1}])$$

Bod, který maximalizuje  $SD(\cdot)$  nazveme nejhlubším simplexovým bodem.

Jinými slovy  $SD(\mathbf{x})$  se rovná pravděpodobnosti, že  $\mathbf{x}$  bude ležet v náhodném simplexu. Je zřejmé, že pro  $\|\mathbf{x}\| \rightarrow \infty$ , kde  $\|\cdot\|$  je Euklidovská norma, platí  $SD(\mathbf{x}) \rightarrow 0$ .

**Definice 4.** Pro  $p$  rozměrný náhodný výběr  $\mathbf{X}_1, \dots, \mathbf{X}_n$  definujeme výběrovou simplexovou hloubku bodu  $\mathbf{x}$  jako:

$$SD_n(\mathbf{x}) = \binom{n}{p+1}^{-1} \sum_{1 \leq i_1 < \dots < i_{p+1} \leq n} I\{\mathbf{x} \in S[\mathbf{X}_{i_1}, \dots, \mathbf{X}_{i_{p+1}}]\}$$

Bod maximalizující  $SD_n(\cdot)$  budeme nazývat výběrovým nejhlubším simplexovým bodem.

Pokud existuje více bodů s maximální hloubkou, budeme brát za nejhlubší bod jejich průměr.

V jednorozměrném případě dostáváme:

$$\begin{aligned} SD(x) &= P(x \in S[X_1, X_2]) \\ &= P(X_1 \leq x \leq X_2) + P(X_2 \leq x \leq X_1) \\ &= 2F(x)(1 - F(x)). \end{aligned}$$

Výraz na pravé straně nabývá maxima pro  $F(x) = 1/2$ , tedy pro  $x$  rovné jednorozměrnému mediánu.

Máme-li body  $\mathbf{x}_1, \dots, \mathbf{x}_{p+1}$ , pak  $\mathbf{x} \in S[\mathbf{x}_1, \dots, \mathbf{x}_{p+1}]$  právě tehdy, když existují jednoznačně určené  $\alpha_1, \dots, \alpha_{p+1}$  takové, že  $\alpha_i \geq 0$ ,  $i = 1, \dots, p+1$ ;  $\alpha_1 + \alpha_2 + \dots + \alpha_{p+1} = 1$ , pro které platí:

$$\mathbf{x} = \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \dots + \alpha_{p+1} \mathbf{x}_{p+1}$$

Odsud ihned dostáváme ekvivarianci vůči afinním transformacím.

Jakou hloubku můžeme očekávat pro simplexový nejhlubší bod? Pro  $p$  rozměrný náhodný vektor  $\mathbf{X}$  s absolutně spojitým úhlově symetrickým rozdělením kolem bodu  $\mathbf{c}$  (tj. náhodné vektory  $(\mathbf{X} - \mathbf{c})/\|\mathbf{X} - \mathbf{c}\|$  a  $-(\mathbf{X} - \mathbf{c})/\|\mathbf{X} - \mathbf{c}\|$  jsou stejně rozdělené), platí  $SD(\mathbf{c}) = 2^{-p}$  a  $SD(\mathbf{x}) \leq 2^{-p} \forall \mathbf{x} \in \mathbb{R}^p$ . Důkaz tohoto tvrzení nalezneme v Liu (1990).

Tamtéž je také uvedeno tvrzení o konzistenci simplexové hloubky. Pokud máme náhodný výběr z rozdělení s omezenou hustotou, pak platí:

$$\sup_{\mathbf{x}} |SD_n(\mathbf{x}) - SD(\mathbf{x})| \xrightarrow{n \rightarrow \infty} 0 \quad s.j.$$

Nabývá-li navíc  $SD(\cdot)$  maxima v právě jednom bodě  $\boldsymbol{\theta}$  a hustota rozdělení je na okolí tohoto bodu nenulová, pak platí:

$$\hat{\boldsymbol{\theta}}_n \xrightarrow{n \rightarrow \infty} \boldsymbol{\theta} \quad s.j.,$$



kde  $\hat{\theta}_n$  značí bod, který maximalizuje  $SD_n(\cdot)$ .

Nevýhodou simplexové hloubky je, že se těžko počítá. Jediný algoritmus, který ho je schopný spočítat v rozumném čase, je algoritmus ISODEPTH (známý též jako AS 307). Ten počítá i halfspace depth zavedenou v části 1.1. Tento algoritmus pracuje jen pro dvojrozměrné vektory.

### 1.3 Mnohorozměrný median (Spatial median)

Mnohorozměrný medián je asi nejpřirozenějším zobecněním jednorozměrného mediánu. Není však robustní a ani ekvivariantní vůči většině afinních transformací.

**Definice 5.** Pro náhodný vektor  $\mathbf{X}$  definujeme mnohorozměrný medián jako

$$\arg \min_{\theta \in \mathbb{R}^p} E \|\mathbf{X} - \theta\|.$$

Pro náhodný výběr  $\mathbf{X}_1, \dots, \mathbf{X}_n$  definujeme výběrový mnohorozměrný medián jako

$$\arg \min_{\theta \in \mathbb{R}^p} \sum_{i=1}^n \|\mathbf{X}_i - \theta\|.$$

$\|\cdot\|$  značí Euklidovskou normu.

Bod zlomu tohoto odhadu je  $1/n$ . Další jeho nevýhodou je, že je ekvivariantní jen vůči afinním transformacím jako je rotace, přetočení os a vynásobení všech složek vektoru stejným číslem. Vůči různé změně měřítka jednotlivých složek náhodného vektoru už ovšem invariantní není. Spočítat jak výběrový, tak i teoretický mnohorozměrný medián je většinou možné jen za pomoci přibližných numerických metod.

Mezi jeho výhody patří vcelku očekávaná vlastnost jednoznačnosti pro spojitá rozdělení (důkaz uveden například v Milasevic & Ducharme (1987)) a fakt, že ze všech mnohorozměrných zobecnění mediánu je mu věnováno asi nejvíce literatury.

V jedné dimenzi je  $\|\cdot\| = |\cdot|$ , takže mnohorozměrný medián se v jednorozměrném případě shoduje s klasickým mediánem.

### 1.4 Další možnosti zavedení mnohorozměrného mediánu

Asi nejjednodušší možností je medián po složkách. Je ho velice jednoduché spočítat, není však ekvivariantní vůči afinním transformacím. Neobstojí ani při rotaci. Je ekvivariantní jen vůči změně měřítka jednotlivých složek náhodného vektoru.

Mnoho dalších možností, založených na různých definicích hloubky dat (*Oja depth*, *Majority depth*, *Convex hull peeling depth*, ...), se dá nalézt v Liu, Jesse & Singh (1999).

## Kapitola 2

# Kvantilová regrese

Při konstrukci a zkoumání vlastností periodických regresních kvantilů použijeme mechanismy kvantilové regrese. V této kapitole by měl být uveden výčet základů a některých pro nás důležitých vlastností kvantilové regrese. Uvedená a mnohá další tvrzení jsou uvedena v Koenker (2005).

V roce 1978 Koenker a Basset rozšířili pojem výběrového kvantilu na regresní model. Jde o vcelku přirozené rozšíření, které v modelu jen s absolutním členem dává výběrový medián.

Mějme regresní model:

$$Y_i = \mathbf{x}_i^T \boldsymbol{\beta} + e_i, \quad i = 1, \dots, n,$$

kde  $Y_1, \dots, Y_n$  jsou odezvy (nebo též náhodná pozorování),  $\boldsymbol{\beta} \in \mathbb{R}^p$ ,  $p < n$  je neznámý parametr,  $\mathbf{x}_i \in \mathbb{R}^p$ ,  $i = 1, \dots, n$  jsou regresory a  $e_1, \dots, e_n$  jsou nezávislé náhodné chyby s rozdělením daným distribuční funkcí  $F$ , pro kterou platí  $F(0) = 1/2$ . Regresní  $\tau$ -kvantil definujeme, pro  $\tau \in (0, 1)$ , jako nějaké řešení minimalizace:

$$\min_{\mathbf{b} \in \mathbb{R}^p} \left[ \tau \sum_{i: Y_i \geq \mathbf{x}_i^T \mathbf{b}} |Y_i - \mathbf{x}_i^T \mathbf{b}| + (1 - \tau) \sum_{i: Y_i < \mathbf{x}_i^T \mathbf{b}} |Y_i - \mathbf{x}_i^T \mathbf{b}| \right]. \quad (2.1)$$

Zavedeme-li funkci  $\rho_\tau$  následujícím způsobem:

$$\rho_\tau(u) = u(\tau - I\{u < 0\}),$$

můžeme minimalizaci v (2.1) přepsat jako:

$$\min_{\mathbf{b} \in \mathbb{R}^p} \sum_{i=1}^n \rho_\tau(Y_i - \mathbf{x}_i^T \mathbf{b}). \quad (2.2)$$

Nechť  $x_{i1} = 1$ ,  $i = 1, \dots, n$ .  $\tau$ -kvantilovou regresní funkcí nazveme (uvažujeme ji jako funkci proměnné  $\mathbf{x}$ ):

$$Q_Y(\tau|\mathbf{x}) = \beta_1 + F^{-1}(\tau) + \beta_1 x_1 + \dots + \beta_p x_p.$$

Vektor  $\hat{\boldsymbol{\beta}}^{(n)}(\tau)$ , který je řešením minimalizace (2.2), bude v tomto případě odhadem vektoru  $(\beta_1 + F^{-1}(\tau), \beta_2, \dots, \beta_p)$ .

Kvantilová regrese a minimalizace (2.2) nám ale umožňuje popis mnohem zajímavějších modelů, než lineárního modelu popsaného na předcházejících řádcích. Model nelineární v parametru  $\beta$ :

$$Y_i = g(\mathbf{x}_i, \beta) + e_i, \quad i = 1, \dots, n,$$

kde  $\beta \in \mathbb{R}^p$ ,  $\mathbf{x}_i \in \mathbb{R}^m$ ,  $g : \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ ,  $e_i$  *i.i.d.* dostaneme minimalizací:

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

Kvantilové regrese se dá ale i použít při odhadech v modelech s nesterjně rozdělenými nezávislými náhodnými chybami. Například na model heteroskedasticity, v kterém rozptyl pozorované veličiny  $Y_i$  závisí na proměnné  $\mathbf{x}_i$ , tvaru

$$Y_i = \mathbf{x}_i^T \beta + \sigma(\mathbf{x}_i, \gamma) e_i, \quad i = 1, \dots, n$$

kde  $e_i$  jsou *i.i.d.*  $F$ ,  $\gamma \in \mathbb{R}^l$  je neznámý parametr a  $\sigma : \mathbb{R}^p \times \mathbb{R}^l \rightarrow \mathbb{R}$  je nějaká nám známá funkce.  $\tau$ -kvantilová funkce má tvar

$$Q_Y(\tau|\mathbf{x}) = \mathbf{x}^T \beta + \sigma(\mathbf{x}, \gamma) F^{-1}(\tau) = (\beta_1 + \sigma(\mathbf{x}, \gamma) F^{-1}(\tau)) + \beta_2 x_2 + \dots + \beta_p x_p.$$

Odhad dostaneme minimalizací

$$\min_{(\mathbf{b}^T, \mathbf{g}^T)^T \in \mathbb{R}^{p+l}} \sum_{i=1}^n \rho_\tau(Y_i - \mathbf{x}_i^T \mathbf{b} - \sigma(\mathbf{x}_i, \mathbf{g})).$$

Označíme-li  $(\hat{\beta}^T, \hat{\gamma}^T)^T$  odhad získaný touto minimalizací, pak odhad  $Q_Y(\tau|\mathbf{x})$  bude vypadat

$$\hat{Q}_Y(\tau|\mathbf{x}) = (\hat{\beta}_1 + \sigma(\mathbf{x}, \hat{\gamma})) + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p.$$

Nakonec uveďme důležitý a zajímavý model s nezávislými a nesterjně rozdělenými náhodnými chybami. Tohoto modelu budeme používat při odhadech periodických regresních kvantilů. Předpokládejme, že máme  $\tau$ -kvantilovou funkci tvaru

$$Q_Y(\tau|\mathbf{x}) = \mathbf{x}^T \beta.$$

To znamená, že např. platí model

$$Y_i = \mathbf{x}_i^T \beta + e_i, \quad i = 1, \dots, n, \tag{2.3}$$

kde  $e_1, \dots, e_n$  jsou nezávislé,  $e_i$  má rozdělení dané distribuční funkcí  $F_i$ , pro kterou  $F_i(0) = \tau$ ,  $i = 1, \dots, n$ . Označme  $\hat{\beta}^{(n)}(\tau)$  vektor řešící minimalizační problém (2.2). Dá se ukázat (viz např. Bantli & Hallin (1999) a Oberhofer (1982)), že za poměrně obecných podmínek je  $\hat{\beta}^{(n)}(\tau)$  konzistentním odhadem parametru  $\beta$ .

Vidíme, že pomocí kvantilové regrese dostaneme konzistentní odhady pro poměrně širokou škálu modelů. Naším cílem však není publikovat výčet modelů v kterých se kvantilová regrese chová „hezky“, ale spíše ukázat, co a jakým způsobem lze v kvantilové regresi odhadovat.

Podívejme se nyní na nějaké vlastnosti odhadu  $p$ -rozměrného parametru  $\beta$  získaným řešením minimalizačního problému (2.2). Předpokládejme, že máme odezvy  $Y_i$ ,  $i = 1, \dots, n$  a jim odpovídající regresory  $\mathbf{x}_i \in \mathbb{R}^p$ ,  $i = 1, \dots, n$ . Problém (2.2) lze převést na úlohu lineárního programování. U těchto úloh je známo, že optimální řešení hledáme v krajních bodech (resp. v lineární kombinaci krajních bodů) množiny omezení. Tato vlastnost implikuje následující větu.

**Věta 6.** Označme  $\hat{\beta} \in \mathbb{R}^p$  optimální řešení úlohy (2.2) a  $\hat{Q}_Y(\tau|\mathbf{x}) = \mathbf{x}^T \hat{\beta}$ . Pak  $Y_i = \hat{Q}_Y(\tau|\mathbf{x}_i)$  pro alespoň  $p$  prvků  $(Y_i, \mathbf{x}_i^T)$ , tzn. že  $\tau$ -kvantilová funkce prochází alespoň  $p$  pozorováními. Pro výběr z absolutně spojitých rozdělení prochází s pravděpodobností jedna právě  $p$  body.

Výběrový  $\tau$ -kvantil v jednorozměrném případě splňuje, že přibližně  $[n\tau]$  pozorování leží pod ním. Platí nějaké podobné tvrzení i pro regresní kvantil? O tom vypovídá následující důležitá věta.

**Věta 7.** Označme  $P$ ,  $N$  a  $Z$  počet kladných, záporných a nulových prvků množiny reziduí  $\{Y_i - \mathbf{x}_i^T \hat{\beta} : i = 1, \dots, n\}$ . Pokud  $x_{i1} = 1$ ,  $i = 1, \dots, n$ , pak

$$N \leq n\tau \leq N + Z$$

a

$$P \leq n(1 - \tau) \leq P + Z.$$

Takže přibližně  $n\tau$  bodů  $(Y_i, \mathbf{x}_i^T)$  leží pod regresní plochou  $\{\hat{Q}_Y(\tau|\mathbf{x}) : \mathbf{x} \in \mathbb{R}^p\}$ .

Od kvantilových funkcí je asi přirozené očekávat, že pro libovolná  $\tau_1, \tau_2$ ,  $0 < \tau_1 < \tau_2 < 1$  bude pro všechna  $\mathbf{x}$  platit:  $Q_Y(\tau_1|\mathbf{x}) < Q_Y(\tau_2|\mathbf{x})$ . A tak by bylo jistě nepříjemným zjištěním, že pro nějaká  $\tau_1, \tau_2$  se odhadnuté regresní funkce  $\hat{Q}_Y(\tau_1|\mathbf{x})$  a  $\hat{Q}_Y(\tau_2|\mathbf{x})$  někde protínají. Tento případ bohužel nelze zcela vyloučit.

Regresní funkce  $Q_Y(\tau|\mathbf{x})$  je pro libovolné  $\mathbf{x}$  neklesající funkcí proměnné  $\tau$ . Podívejme se pro jaké hodnoty  $\mathbf{x}$  lze tuto vlastnost zaručit i pro její odhad.

**Věta 8.** Označme  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ . Pak  $\hat{Q}_Y(\tau|\bar{\mathbf{x}})$  je neklesající funkcí v proměnné  $\tau \in (0, 1)$ .

Tvrzení věty nám samozřejmě nezaručí, že  $\hat{Q}_Y(\tau|\mathbf{x})$  bude neklesající v  $\tau$  pro všechna  $\mathbf{x}$ . Ale je vidět, že pokud bude tato vlastnost porušena, pak to pro „rozumně“ rozmístěná data nastane poměrně daleko od  $\bar{\mathbf{x}}$ .

Ještě uvedme větu o ekvivarianci odhadu kvantilové regrese. Označme  $\hat{\beta}(\tau; \mathbf{Y}, \mathbf{X})$  odhad založený na pozorováních  $(\mathbf{Y}, \mathbf{X})$ ,  $\mathbf{X}$  je matice jejíž řádky odpovídají regresorům  $\mathbf{x}_i$  a  $\mathbf{Y}$  je vektor pozorování.

**Věta 9.** Nechť  $A \in \mathbb{R}^{p \times p}$  je regulární matice,  $\gamma \in \mathbb{R}^p$ ,  $a > 0$  a  $h$  neklesající funkce. Pak pro libovolné  $\tau \in (0, 1)$  platí

$$(i) \hat{\beta}(\tau; a\mathbf{Y}, \mathbf{X}) = a\hat{\beta}(\tau; \mathbf{Y}, \mathbf{X})$$

$$(ii) \hat{\beta}(\tau; -a\mathbf{Y}, \mathbf{X}) = -a\hat{\beta}(\tau; \mathbf{Y}, \mathbf{X})$$

$$(iii) \hat{\beta}(\tau; \mathbf{Y} + \mathbf{X}\gamma, \mathbf{X}) = \hat{\beta}(\tau; \mathbf{Y}, \mathbf{X}) + \gamma$$

$$(iv) \hat{\beta}(\tau; \mathbf{Y}, \mathbf{X}A) = A^{-1}\hat{\beta}(\tau; \mathbf{Y}, \mathbf{X})$$

$$(v) Q_{h(Y)}(\tau|\mathbf{x}) = h(Q_Y(\tau|\mathbf{x}))$$

Narozdíl od odhadu metodou nejmenších čtverců, který je dost citlivý na odlehlá pozorování  $Y_i$ , je odhad kvantilové regrese velice robustní. Při posunutí jakéhokoliv pozorování  $Y_i$  libovolně daleko od regresní plochy se odhad nezmění. Tuto vlastnost přesně popisuje následující věta.

**Věta 10.** *Pro libovolnou diagonální matici  $D$ , která má na diagonále nezáporné prvky, platí*

$$\hat{\beta}(\tau; \mathbf{Y}, \mathbf{X}) = \hat{\beta}(\tau; \mathbf{X}) + D\hat{\mathbf{u}},$$

kde  $\hat{\mathbf{u}} = \mathbf{Y} - \mathbf{X}\hat{\beta}(\tau; \mathbf{Y}, \mathbf{X})$  je vektor odhadu reziduí.

K podobnému závěru, co se robustnosti týká, dojdeme i výpočtem influenční funkce. Ta je omezená ve veličině odpovídající odezvě  $Y_i$ . A to také značí malou citlivost na odlehlé hodnoty  $Y_i$ . Na druhou stranu není omezená ve veličině odpovídající regresorům  $\mathbf{x}_i$  a tedy odlehlé hodnoty  $\mathbf{x}_i$  mohou mít na odhad někdy nezanedbatelný vliv.

## Kapitola 3

# Periodické regresní kvantily

V této kapitole se budeme věnovat novému způsobu konstrukce konfidenčních množin na dané hladině. Jak její konstrukci pro náhodné veličiny se spojitým rozdělením, tak pro její výběrovou variantu, a na jejich vzájemný vztah. Postup je založen na transformaci vycentrovaných dat do polárních souřadnic a poté určení jakýchsi směrových kvantilů. V teoretickém případě to provedeme pomocí rozdělení podmíněného volbou směru (resp. volbou úhlů, který tento směr určí) z daného centrálního bodu, ve výběrové variantě pomocí kvantilové regrese a použití trigonometrických řad. Pro vycentrování dat a tedy určení bodu, který by byl pro nás v nějakém smyslu středem našich dat, použijeme nějaký z parametrů polohy popsanych v kapitole 1. Pro výběr tohoto středu se mi zdál nejvhodnější *nejhlubší bod*. Bude proto v postupech popsanych v následujícím textu použit tento bod.

### 3.1 Teoretické periodické regresní kvantily

Mějme  $k$ -rozměrný ( $k > 1$ ) náhodný vektor  $\mathbf{X}$  se spojitým rozdělením daným hustotou  $f(\mathbf{x})$ . Označme  $\boldsymbol{\theta}$  teoretický nejhlubší bod tohoto rozdělení. Pro každou polopřímku s počátkem v bodě  $\boldsymbol{\theta}$  najdeme podmíněné rozdělení (podmíněné volbou této polopřímky) náhodné veličiny udávající vzdálenost bodů ležících na této polopřímce od počátku. Pokusme se nalézt takovou hodnotu  $r$  pro kterou bude pravděpodobnost, že vzdálenost bodů od počátku je menší než  $r$ , rovna nějaké zvolené hodnotě  $\tau \in (0, 1)$ . To provedeme pomocí transformace do *polárních souřadnic* (v prostorech vyšších dimenzí se používá též označení *hypersférické souřadnice*). Každý bod bude charakterizován veličinou  $\rho$  udávající vzdálenost od počátku a úhly  $\phi_1, \dots, \phi_{k-1}$ , které udávají směr, kterým se z počátku vydáváme. Tyto veličiny jsou určeny vztahy:

$$\begin{aligned} X_1 &= \theta_1 + \rho \sin \phi_1 \sin \phi_2 \cdots \sin \phi_{k-2} \sin \phi_{k-1}, \\ X_2 &= \theta_2 + \rho \sin \phi_1 \sin \phi_2 \cdots \sin \phi_{k-2} \cos \phi_{k-1}, \\ &\vdots \\ X_{k-1} &= \theta_{k-1} + \rho \sin \phi_1 \cos \phi_2, \\ X_k &= \theta_k + \rho \cos \phi_1, \end{aligned}$$

kde

$$\phi_i \in (0, \pi), \quad i = 1, \dots, k-2, \quad \phi_{k-1} \in [0, 2\pi) \text{ a } \rho > 0.$$

Jakobián je roven

$$J = \rho^{k-1} \sin^{k-2} \phi_1 \cdots \sin \phi_{k-2}.$$

Hustota vektoru  $(\rho, \phi_1, \dots, \phi_{k-1})^T = (\rho, \phi^T)^T$  bude tedy rovna

$$\begin{aligned} p(r, \varphi_1, \dots, \varphi_{k-1}) &= \\ &= r^{k-1} |\sin^{k-2} \varphi_1 \cdots \sin \varphi_{k-2}| f(\theta_1 + r \sin \varphi_1 \sin \varphi_2 \cdots \sin \varphi_{k-2} \sin \varphi_{k-1}, \dots, \theta_k + r \cos \varphi_1). \end{aligned}$$

Marginální hustotu náhodného vektoru  $\phi$  je

$$s(\varphi) = \int_0^{+\infty} p(r, \varphi) dr.$$

A pro podmíněnou hustotu náhodné veličiny  $\rho$  při daném  $\phi = \varphi$  tedy platí

$$q(r|\varphi) = \begin{cases} \frac{p(r, \varphi)}{s(\varphi)} & \text{pro } s(\varphi) \neq 0, \\ 0 & \text{pro } s(\varphi) = 0. \end{cases}$$

Pro dané  $\varphi$  a dané  $\tau$  hledáme hodnotu  $r(\tau|\varphi)$  pro kterou platí

$$\tau = P(\rho \leq r(\tau|\varphi) | \phi = \varphi) = \int_0^{r(\tau|\varphi)} q(r|\varphi) dr = Q(r(\tau|\varphi)|\varphi),$$

kde  $Q(\cdot|\varphi)$  značí distribuční funkci veličiny  $\rho$  při  $\phi = \varphi$ . Označíme-li  $Q^{-1}(\cdot|\varphi)$  inverzi této distribuční funkce, pak definujeme *směrový  $\tau$ -kvantil* pro dané  $\varphi$  jako

$$r(\tau|\varphi) = Q^{-1}(\tau|\varphi).$$

Nyní můžeme přejít zpět ke kartézským souřadnicím.

**Definice 6.** Označme  $\mathcal{M} = \{\mathbf{x} : f(\mathbf{x}) > 0\}$ . Teoretickým periodickým regresním kvantilem nazveme množinu

$$\mathcal{K}(\tau) = \mathcal{C}(\tau) \cap \mathcal{M},$$

kde

$$\begin{aligned} \mathcal{C}(\tau) = \{\mathbf{x} \in \mathbb{R}^k : \begin{aligned} x_1 &= \theta_1 + r \sin \varphi_1 \sin \varphi_2 \cdots \sin \varphi_{k-2} \sin \varphi_{k-1}, \\ x_2 &= \theta_2 + r \sin \varphi_1 \sin \varphi_2 \cdots \sin \varphi_{k-2} \cos \varphi_{k-1}, \\ &\vdots \\ x_{k-1} &= \theta_{k-1} + r \sin \varphi_1 \cos \varphi_2, \\ x_k &= \theta_k + r \cos \varphi_1, \\ 0 &\leq \varphi_1, \dots, \varphi_{k-2} \leq \pi, \quad 0 \leq \varphi_{k-1} < 2\pi, \quad 0 \leq r \leq r(\tau|\varphi) \end{aligned}\} \end{aligned} \tag{3.1}$$

Má množina  $\mathcal{K}(\tau)$  vlastnosti, které bychom od konfidenční množiny mohli očekávat? Zkusme se podívat na některé její vlastnosti.

**Věta 11.**

$$P(\mathbf{X} \in \mathcal{K}(\tau)) = \tau.$$

*Důkaz.*

$$\begin{aligned} P(\mathbf{X} \in \mathcal{K}(\tau)) &= P(0 \leq \phi_1, \dots, \phi_{k-2} \leq \pi, 0 \leq \phi_{k-1} < 2\pi, 0 \leq \rho \leq r(\tau|\phi)) \\ &= \int_0^{2\pi} \int_0^\pi \cdots \int_0^\pi \int_0^{r(\tau|\varphi)} q(r|\varphi) s(\varphi) dr d\varphi_1 \dots d\varphi_{k-1} \\ &= \tau \int_0^{2\pi} \int_0^\pi \cdots \int_0^\pi s(\varphi) d\varphi_1 \dots d\varphi_{k-1} \\ &= \tau. \end{aligned}$$

□

**Věta 12.** *Množina  $\mathcal{K}(\tau)$  je omezená.*

*Důkaz.* Okamžitě plyne z toho, že pro spojitá rozdělení platí  $r(\tau|\varphi) < +\infty$  s.j.  $\forall \varphi$ . □

**Věta 13.** *Pokud  $0 < \tau_1 < \tau_2 < 1$ , pak  $\mathcal{K}(\tau_1) \subset \mathcal{K}(\tau_2)$ .*

*Důkaz.* Zřejmé. □

Povrchem množiny  $\mathcal{K}(\tau)$  nazveme plochu, která vznikne nahradíme-li ve (3.1) hodnotu  $r$  směrovým kvantilem  $r(\tau|\varphi)$ . Od povrchu konfidenčních množin je asi přirozené požadovat, aby byl spojitý, tj. aby funkce  $r(\tau|\varphi)$  byla spojitou funkcí proměnné  $\varphi$  pro všechna  $\tau \in (0, 1)$ . Navíc je ještě důležitý tvar množiny  $\mathcal{M}$  a poloha bodu  $\theta$  vůči této množině. Směrové kvantily hledáme na polopřímkách začínajících v bodě  $\theta$ . Proto se omezíme na tzv. hvězdicovité množiny. Jak bude vidět v příkladu 4 pro jiné množiny nemusí množina  $\mathcal{K}(\tau)$  mít zrovna uspokojivý tvar.

**Definice 7.** *Množinu  $\mathcal{S} \subset \mathbb{R}^k$  nazveme hvězdicovitou kolem  $\xi$ , pokud pro každé  $x \in \mathcal{S}$ , úsečka spojující body  $\xi$  a  $x$  je celá obsažena v  $\mathcal{S}$ .*

A z konstrukce množiny  $\mathcal{K}(\tau)$  je zřejmé následující jednoduché tvrzení.

**Věta 14.** *Pokud je  $\mathcal{M}$  hvězdicovitá, pak je i  $\mathcal{K}(\tau)$  hvězdicovitá.*

Je zřejmé, že pro hvězdicovité množiny  $\mathcal{M}$  může mít funkce  $r(\tau|\varphi)$  několik bodů nespojitosti. Pokusme se situaci ilustrovat na následujících třech příkladech.

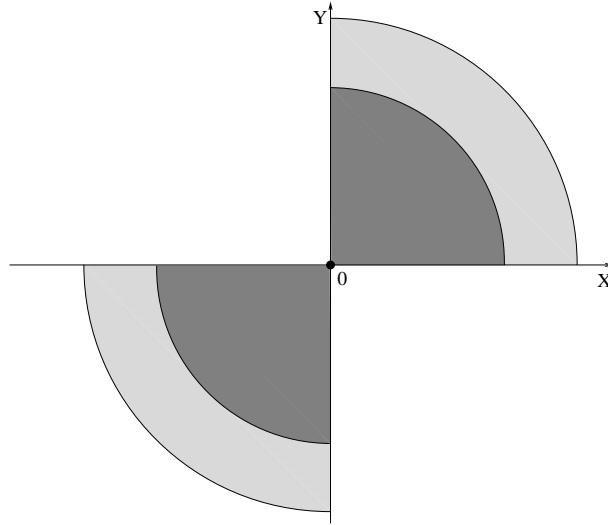
### Příklad 3

Představme si vektor  $(X, Y)^T$  s rovnoměrným rozdělením na dvou protilehlých kruhových výsečích z jednotkového kruhu se středem v počátku, tak jak je to zobrazeno na obr. 3.1. Nejhlubší bod se shoduje s počátkem. Množina  $\mathcal{M}$  je hvězdicovitá kolem počátku.

Hustota vektoru  $(X, Y)^T$  přetransformovaného do polárních souřadnic je rovna

$$p(r, \varphi) = \frac{r}{2\pi} I\{r^2 \leq 1, \varphi \in [0, \pi/2] \cup [\pi, 3\pi/2]\}.$$





Obrázek 3.1: Teoretický periodický regresní kvantil pro rovnoměrné rozdělení na kruhových výsečích.

Pro podmíněnou hustotu a distribuční funkci platí

$$\begin{aligned} q(r|\varphi) &= \frac{r}{\frac{2\pi}{1}} I\{r^2 \leq 1, \varphi \in [0, \pi/2] \cup [\pi, 3\pi/2]\} \\ &= 2r I\{r^2 \leq 1, \varphi \in [0, \pi/2] \cup [\pi, 3\pi/2]\}, \\ Q(r|\varphi) &= r^2 I\{r^2 \leq 1, \varphi \in [0, \pi/2] \cup [\pi, 3\pi/2]\}. \end{aligned}$$

Z toho nakonec dostáváme  $\tau$ -směrový kvantil ve tvaru

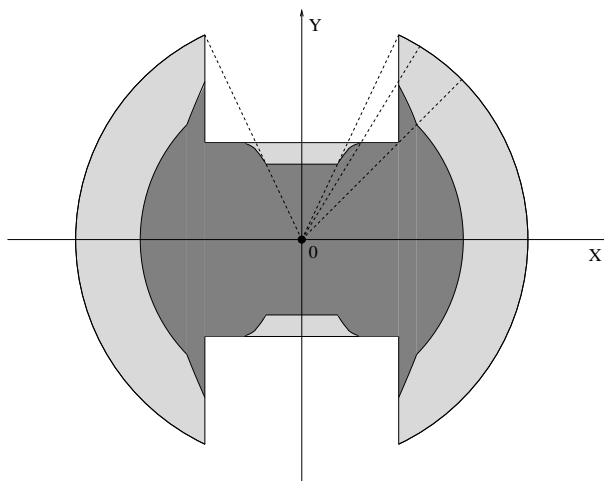
$$r(\tau|\varphi) = \sqrt{\tau} I\{\varphi \in [0, \pi/2] \cup [\pi, 3\pi/2]\}.$$

$\mathcal{K}(\tau)$  se bude skládat ze dvou kruhových výsečí o poloměru  $\sqrt{\tau}$ . Na obr. 3.1 je tmavě šedou barvou vybarvena množina  $\mathcal{K}(0.5)$ .  $r(\tau|\varphi)$  je spojitá na intervalech  $(0, \pi/2)$  a  $(\pi, 3\pi/2)$ . Mimo tyto intervaly není definována, protože sdružená hustota  $p(r, \varphi)$  je zde nulová.

#### Příklad 4

Mějme znovu rovnoměrné rozdělení, ale teď na množině, která nemá hvězdicovitý tvar. Množina  $\mathcal{M}$  vznikla vyříznutím dvou obdélníků z kruhu. Strana tohoto obdélníku rovnoběžná s osou  $X$  je od ní vzdálena v poměru  $5/11$  poloměru a má velikost  $10/11$  poloměru kružnice. Situace je zobrazena na obr. 3.2. Po delším počítání se dobereme k množině  $\mathcal{K}(0.5)$  tvaru, který je na obrázku vyznačen tmavě šedou barvou. Oproti předchozímu příkladu už tvar zkonstruované konfidenční množiny nemusí být zrovna pro leckoho přijatelný.  $r(\frac{1}{2}|\varphi)$  je zde definována pro  $\varphi \in [0, 2\pi)$ , ale má několik bodů nespojitosti.

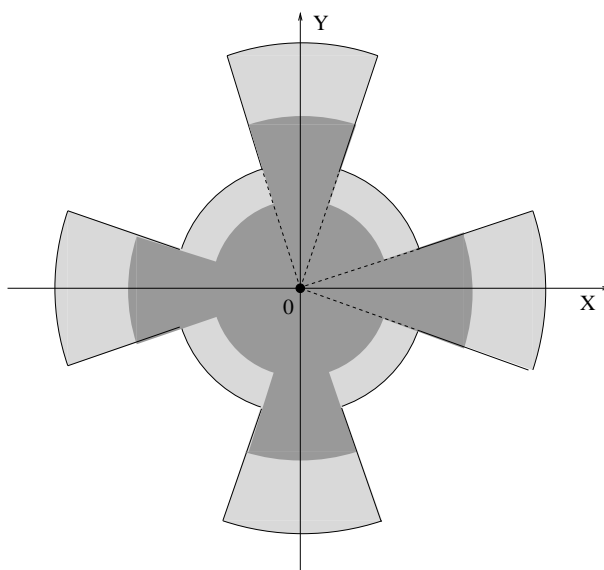
Následující příklad ilustruje, že ani hvězdicovitý tvar množiny  $\mathcal{M}$  spolu s vlastností, že bod  $\theta$  je vnitřním bodem množiny  $\mathcal{M}$ , nezaručuje spojitost směrového kvantilu  $r(\tau|\cdot)$  a ani obzvlášť pěkný tvar množiny  $\mathcal{K}(\tau)$ .



Obrázek 3.2: Teoretický periodický regresní kvantil pro rovnoměrné rozdělení na množině nehvězdicovitého tvaru.

#### Příklad 5

Pro rovnoměrné rozdělení na množině tvaru větráku - viz obr. 3.3 - dostaneme směrový kvantil podobnými výpočty jako v příkladu 3. V tomto případě je funkce  $r(\tau|\varphi)$  definována



Obrázek 3.3: Nespojitý teoretický periodický regresní kvantil pro rovnoměrné rozdělení na hvězdicovité množině.

pro všechna  $\varphi \in [0, 2\pi)$ , v „listech“ větráku bude rovna  $r_1\sqrt{\tau}$  a ve „středu“ větráku  $r_2\sqrt{\tau}$ .  $r_1$  značí délku listu,  $r_2$  poloměr středu větráku. Je tedy spojitá mimo konečně mnoho bodů. Tyto body jsme dostali díky tvaru množiny  $\mathcal{M}$ , kde listy větráku jsou vlastně kruhové výseče se středem v počátku. Takže bod nespojitosti vlastně určuje polopřímku vymežující okraj těchto výsečí. Na jedné straně od tohoto bodu hledáme směrový kvantil podobně jako v kruhu

o poloměru  $r_1$ , na druhé straně jako v kruhu o poloměru  $r_2$ . Pro podobnou situaci, kdy by listy větráku nebyly kruhovými výsečemi se středem v počátku, ale např. obdélníky, bychom už dostali směrové kvantily spojitě (v takové situaci by totiž polopřímka začínající v počátku protínala okraj listu právě v jednom bodě). Tmavě šedá plocha na obr. 3.3 opět vymezuje množinu  $\mathcal{K}(0.5)$ .

V předcházejících příkladech jsme mimo jiné viděli, že spojitost funkce  $r(\tau|\cdot)$  závisí nejen na typu rozdělení náhodného vektoru, ale také na tvaru množiny  $\mathcal{M}$  a poloze bodu  $\theta$  vůči této množině.

**Věta 15** (Spojitost). *Nechť je množina  $\mathcal{M}$  hvězdicovitá kolem bodu  $\theta$  a souvislá. Nechť je  $f$  na ní spojitá. Označme  $\partial\mathcal{M}$  hranici  $\mathcal{M}$ . Pokud pro každou přímku  $l$  procházející bodem  $\theta$  platí, že množina  $l \cap \partial\mathcal{M}$  má jen konečně mnoho bodů, pak je povrch množiny  $\mathcal{K}(\tau)$  spojitý (tj. funkce  $r(\tau|\cdot)$  je spojitá) pro všechna  $\tau \in (0, 1)$ .*

*Důkaz.* Stačí ověřit, že pro libovolnou pevně zvolenou hodnotu  $u$  je

$$G : \varphi \mapsto \int_0^u q(r|\varphi) dr = Q(u|\varphi)$$

spojitou funkcí. Pak již se dá ukázat, že je spojitá i funkce  $Q^{-1}(\tau|\varphi) = r(\tau|\varphi)$  (jako funkce proměnné  $\varphi$ ) pro libovolnou pevně zvolenou hodnotu  $\tau$ .

Pro spojitost  $G$  potřebujeme ověřit integrovatelnost (absolutní)  $q(\cdot|\varphi)$  pro všechna  $\varphi$  (zřejmé), měřitelnost  $q(\cdot|\varphi)$  pro s.v.  $\varphi$  (téměř zřejmé) a spojitost  $q(r|\cdot)$  pro s.v.  $r$ . Díky předpokladům věty platí pro libovolné  $\varphi, \tilde{\varphi}$

$$|q(r|\varphi) - q(r|\tilde{\varphi})| = \left| \frac{p(r, \varphi)}{s(\varphi)} - \frac{p(r, \tilde{\varphi})}{s(\tilde{\varphi})} \right| \leq \text{konst } r^{k-1} |f(\mathbf{x}) - f(\tilde{\mathbf{x}})|.$$

Spojitost  $q(r|\cdot)$  tedy plyne ze spojitosti  $f$ . □

Ještě uvedme důležitou větu o ekvivarianci.

**Věta 16** (Ekvivariance). *Periodický regresní kvantil je ekvivariantní vzhledem k afinním transformacím.*

*Důkaz.* Ekvivariance vzhledem k posunutí, rotaci a překlopení je zřejmá.

Libovolnou afinní transformací se zobrazuje přímka na přímku. Navíc se zachovává rovnoběžnost přímk, rovin, atd. Díky tomu bude po obecné afinní transformaci podmíněné rozdělení  $\rho$  při pevné  $\phi$  oproti původnímu jen natažené resp. zúžené (ve smyslu vynásobení konstantou  $a > 1$  resp.  $0 < a < 1$ ). Odtud a z konstrukce periodických kvantilů už je ekvivariance viditelná. □

### Příklad 6 Exponenciální rozdělení

Vraťme se k příkladu 2. Souřadnice nejhlubšího bodu jsou rovny  $\theta_0 \doteq 0.76307$ . Hustota vektoru  $(X, Y)^T$  je  $f(x, y) = e^{-x-y} I\{x > 0, y > 0\}$ . Náhodné veličiny  $\rho$  a  $\phi$  dostaneme z přechodu k polárním souřadnicím

$$\begin{aligned} X &= \theta_0 + \rho \cos \phi, \\ Y &= \theta_0 + \rho \sin \phi. \end{aligned}$$

Jejich sdružená hustota je

$$p(r, \varphi) = r e^{-2\theta_0} e^{-r(\cos \varphi + \sin \varphi)} I\{\theta_0 + r \cos \varphi > 0, \theta_0 + r \sin \varphi > 0\}.$$

Hustota  $p(r, \varphi)$  bude nenulová pro hodnoty  $\varphi, r$  takové, že:

1. Pokud  $\varphi \in [0, \pi/2)$ , pak  $r > 0$ .
2. Pro  $\varphi \in [\pi/2, \frac{5}{4}\pi)$  bude  $0 < r < \frac{\theta_0}{\cos(\varphi - \pi)}$ .
3. Pro  $\varphi \in [\frac{5}{4}\pi, 2\pi)$  bude  $0 < r < \frac{\theta_0}{\cos(\varphi - \frac{3}{2}\pi)}$ .

Postupnými výpočty dojdeme k podmíněné distribuční funkci  $\rho$  při daném  $\phi = \varphi$  tvaru:

$$Q(r|\varphi) = \begin{cases} H(r|\varphi) & \text{pro } \varphi \in [0, \pi/2) \\ \frac{H(r|\varphi)}{1 - [1 + \frac{\theta_0}{\cos(\varphi - \pi)}(\cos \varphi + \sin \varphi)] \exp\{-\frac{\theta_0}{\cos(\varphi - \pi)}(\cos \varphi + \sin \varphi)\}} & \text{pro } \varphi \in [\pi/2, \frac{5}{4}\pi) \\ \frac{H(r|\varphi)}{1 - [1 + \frac{\theta_0}{\cos(\varphi - \frac{3}{2}\pi)}(\cos \varphi + \sin \varphi)] \exp\{-\frac{\theta_0}{\cos(\varphi - \frac{3}{2}\pi)}(\cos \varphi + \sin \varphi)\}} & \text{pro } \varphi \in [\frac{5}{4}\pi, 2\pi) \end{cases}$$

kde  $H(r|\varphi) = 1 - [1 + r(\cos \varphi + \sin \varphi)] \exp\{-r(\cos \varphi + \sin \varphi)\}$ . Inverzní funkci k  $H(r|\varphi)$  bohužel nelze vyjádřit explicitně. A tak směrový kvantil  $r(\tau|\varphi)$  budeme hledat pro pevnou hodnotu  $\varphi$  jako numerické řešení rovnice

$$Q(r|\varphi) = \tau.$$

Exponenciální rozdělení splňuje předpoklady věty o spojitosti, a tak dostáváme, že křivka, kterou získáme pomocí  $r(\tau|\varphi)$  přechodem zpět ke kartézským souřadnicím, bude spojitá a uzavřená. Vypočtené konfidenční množiny  $\mathcal{K}(\tau)$ ,  $\tau = 0.1, 0.2, \dots, 0.9$  a nejhlubší bod jsou zobrazeny na obr. 3.4.

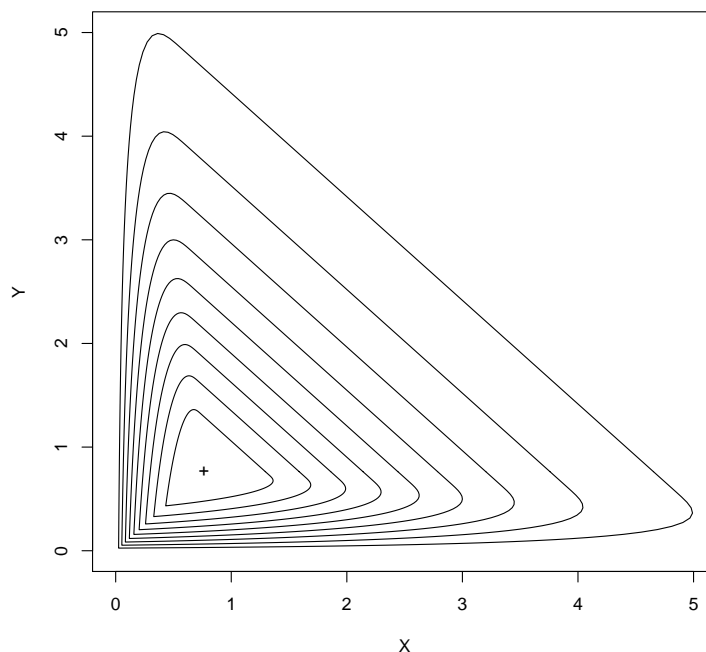
## 3.2 Výběrové periodické regresní kvantily

Funkce  $r(\tau|\varphi)$  zavedená v předchozí kapitole je v dvojrozměrném případě  $2\pi$  periodickou funkcí v proměnné  $\varphi$ . Za dosti obecných předpokladů existuje Fourierův rozvoj  $r(\tau|\varphi)$ , který k této funkci konverguje stejnoměrně. Pro náhodný výběr se tedy nabízí použít k jejímu odhadu trigonometrické řady s konečným počtem sčítanců. Koeficienty této řady se budeme snažit získat pomocí kvantilové regrese.

Ve vícerozměrném případě je situace podobná. V textu se zaměříme hlavně na trojrozměrný případ, kde se nejprve seznámíme se základy Fourierových rozvojevů funkcí dvou proměnných. Oproti dvojrozměrnému případu budou na koeficienty trigonometrické řady, kterou použijeme k odhadu, kladeny další požadavky tak, aby získaný odhad splňoval podobné vlastnosti jako jeho teoretický protějšek. Princip konstrukce trojrozměrných výběrových periodických kvantilů lze analogicky provést i pro více než trojrozměrné náhodné výběry.

### Dvojrozměrný případ

Ještě než začneme popisovat konstrukci dvojrozměrných periodických kvantilů, připomeňme si nějaké základní pojmy teorie Fourierových řad jedné proměnné.



Obrázek 3.4: Teoretické periodické regresní kvantily pro exponenciální rozdělení,  $\tau = 0.1, 0.2, \dots, 0.9$

Mějme  $2\pi$ -periodickou funkci  $f : \mathbb{R} \rightarrow \mathbb{R}$ , která je lebesgueovsky integrovatelná na omezených intervalech. Fourierovými koeficienty funkce  $f$  budeme rozumět čísla

$$\alpha_n = \frac{1}{\pi} \int_0^{2\pi} f(x) \cos nx \, dx, \quad n = 0, 1, \dots,$$

$$\beta_n = \frac{1}{\pi} \int_0^{2\pi} f(x) \sin nx \, dx, \quad n = 1, 2, \dots$$

Fourierovou řadou funkce  $f$  nazveme řadu

$$\frac{\alpha_0}{2} + \sum_{n=1}^{\infty} (\alpha_n \cos nx + \beta_n \sin nx). \quad (3.2)$$

Uvedme ještě poměrně silnou, ale pro další postup plně vyhovující, postačující podmínku pro konvergenci Fourierovy řady.

**Věta 17.** *Nechť  $f$  splňuje podmínky popsané v úvodu této sekce a nechť je navíc spojitá a po částech hladká (v tom smyslu, že existuje jen konečně mnoho bodů v kterých neexistuje derivace funkce). Pak Fourierova řada funkce  $f$  konverguje stejnoměrně k funkci  $f$  na  $\mathbb{R}$ .*

**Poznámka.** *V části 3.1 jsme uvedli příklady nespojitých po částech hladkých směrových kvantilů. Dá se ukázat, že pro tyto funkce konverguje jejich řada bodově k funkci  $g$ , kde*

$$g(x) = \frac{1}{2} \lim_{h \rightarrow 0^+} (f(x+h) + f(x-h)).$$

Pro většinu běžně používaných spojitých rozdělení jsou předpoklady spojitosti a částečné hladkosti splněny.

Fourierova řada směrového kvantilu  $r(\tau|\varphi)$  tedy za těchto předpokladů k němu stejnoměrně konverguje. Tuto řadu se pokusíme odhadnout trigonometrickou řadou s konečným počtem sčítanců. Pro nějaké  $p \in \mathbb{N}, p < \infty$  a  $\tau \in (0, 1)$  se tedy budeme snažit získat odhad koeficientů v řadě

$$r_p(\tau|\varphi) = a_0 + \sum_{j=1}^p (a_j \cos j\varphi + b_j \sin j\varphi), \quad (3.3)$$

tak aby byla co nejvíce podobná funkci  $r(\tau|\varphi)$ .

Je ihned vidět, že funkce  $r_p$  splňuje požadavky, které bychom od odhadu funkce  $r(\tau|\varphi)$  přirozeně očekávali - spojitost a  $2\pi$  periodicitu. Koeficienty v (3.3) se pokusíme odhadnout pomocí kvantilové regrese, jejíž základy jsou popsány v kapitole 2.

Mějme dvojrozměrný náhodný výběr  $(X_i, Y_i)^T$ ,  $i = 1, \dots, n$ . Výběrový nejhlubší bod tohoto výběru označme  $(\hat{\theta}_1, \hat{\theta}_2)^T$ . Přejdeme k dvojrozměrnému vektoru  $(R_i, F_i)^T$ ,  $i = 1, \dots, n$  splňujícímu

$$\begin{aligned} X_i &= \hat{\theta}_1 + R_i \cos F_i, \quad i = 1, \dots, n, \\ Y_i &= \hat{\theta}_2 + R_i \sin F_i, \quad i = 1, \dots, n \end{aligned}$$

a  $R_i > 0$ ,  $F_i \in [0, 2\pi)$ ,  $i = 1, \dots, n$ .

Nyní můžeme pomocí kvantilové regrese odhadnout  $\tau$ -regresní kvantil pro odezvu a regresory tvaru

$$\begin{pmatrix} R_1 \\ R_2 \\ \vdots \\ R_n \end{pmatrix}, \quad \begin{pmatrix} 1 & \cos F_1 & \cos 2F_1 & \dots & \cos pF_1 & \sin F_1 & \sin 2F_1 & \dots & \sin pF_1 \\ 1 & \cos F_2 & \cos 2F_2 & \dots & \cos pF_2 & \sin F_2 & \sin 2F_2 & \dots & \sin pF_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \cos F_n & \cos 2F_n & \dots & \cos pF_n & \sin F_n & \sin 2F_n & \dots & \sin pF_n \end{pmatrix} \quad (3.4)$$

Postupujeme vlastně obdobně jako při odhadu koeficientů v modelu s polynomickou závislostí na jednom regresoru:  $Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \dots + \beta_p X_i^p + e_i$ ,  $i = 1, \dots, n$ . Jen místo  $x^k$  použijeme funkce  $\cos kx$  a  $\sin kx$ .

Je zřejmé, že pro  $p$  musí platit podmínka  $2p + 1 < n$ .

**Definice 8.** Označme  $\hat{\beta} = (\hat{\alpha}_0, \hat{\alpha}_1, \dots, \hat{\alpha}_p, \hat{\beta}_1, \dots, \hat{\beta}_p)^T \in \mathbb{R}^{2p+1}$  odhad minimalizační úlohy (2.2) (kvantilová regrese) pro regresory a odezvu dané v (3.4) a dané  $\tau$ . Výběrovým směrovým  $\tau$ -kvantilem řádu  $p$  nazveme funkci

$$r_p(\tau|\varphi) = \hat{\alpha}_0 + \sum_{j=1}^p \hat{\alpha}_j \cos j\varphi + \hat{\beta}_j \sin j\varphi.$$

Množinu

$$\begin{aligned} \mathcal{K}_p(\tau) = \{(x, y)^T \in \mathbb{R}^2 : & x = \hat{\theta}_1 + r \cos \varphi, \\ & y = \hat{\theta}_2 + r \sin \varphi, \\ & 0 \leq \varphi < 2\pi, \quad 0 \leq r \leq r_p(\tau|\varphi)\} \end{aligned}$$

nazveme výběrovým periodickým regresním kvantilem.

Jak již bylo řečeno, tak řád odhadu  $p$  periodických regresních kvantilů je omezen shora číslem  $\frac{n-1}{2}$ . Ve skutečnosti se budeme snažit volit řád daleko nižší. Věta 6 říká, že pro výběry ze spojitých rozdělání bude funkce  $r_p(\tau|\varphi)$  procházet  $2p+1$  body. Čím větší je toto číslo, tím více body funkce  $r_p$  prochází. A je zřejmé, že s rostoucím počtem těchto bodů se bude  $r_p$  více a více vlnit. Pro krajní případ  $p = \lfloor \frac{n-1}{2} \rfloor$  dostaneme velice „divokou“ funkci, která prochází všemi pozorovanými body. Na druhou stranu pro  $p$  velmi malé (vzhledem k  $n$ ) hrozí, že  $r_p$  bude příliš hrubým odhadem  $r$ . Určit pro dané  $n$  velikost řádu  $p$  není jednoduché. Zatím si řekneme, že  $p$  budeme volit výrazně menší než  $n$ . Poznamenejme, že už pro  $p = 10$  je množina všech funkcí daných v (3.3) poměrně bohatá.

Dále budeme požadovat, aby platilo

$$p \xrightarrow{n \rightarrow \infty} \infty \quad \text{a} \quad \frac{p}{n^s} \xrightarrow{n \rightarrow \infty} 0 \quad (3.5)$$

pro nějaké  $0 < s \leq 1$ .

Následující tvrzení okamžitě vyplývají z konstrukce výběrových kvantilů.

**Věta 18.** *Funkce  $r_p(\tau|\cdot)$  je spojitá a  $2\pi$  periodická. Množina  $\mathcal{K}_p(\tau)$  je hvězdicovitá a omezená.*

Z věty 7 dostáváme následující důležitou větu o počtu bodů ležících v  $\mathcal{K}_p(\tau)$ .

**Věta 19.** *Označme  $N$  počet bodů  $(X_i, Y_i)$ ,  $i = 1, \dots, n$  ležících ve vnitřku množiny  $\mathcal{K}_p(\tau)$  a  $Z$  počet bodů ležících na její hranici. Pak*

$$n\tau - Z \leq N \leq n\tau.$$

Pro výběry ze spojitých rozdělání je  $Z = 2p + 1$ . Pokud navíc budeme volit velikost řádu  $p$  tak, aby splňovala (3.5) pak

$$\frac{N}{n} \xrightarrow{n \rightarrow \infty} \tau.$$

Posunutí, rotace, stejná změna měřítka všech složek náhodného vektoru (zvětšení, zmenšení) a překlopení kolem libovolné přímky patří mezi transformace, které nijak výrazně nemění rozvržení dat. A tak by asi bylo přirozené od dobrých odhadů očekávat, aby byly vůči těmto transformacím ekvivariantní.

Ekvivariance vzhledem k posunutí je zřejmá.

U stejné změny měřítka jednotlivých složek je situace také jednoduchá. Pro  $k > 0$  dostáváme náhodný výběr  $(kX_i, kY_i)$ ,  $i = 1, \dots, n$ . Výběrový nejhlubší bod bude  $(k\hat{\theta}_1, k\hat{\theta}_2)^T$ . Úhel charakterizující jednotlivá pozorování zůstává stejný jako v původním výběru. Vzdálenost  $(kX_i, kY_i)$  od středu je rovna  $\sqrt{k^2((X_i - \hat{\theta}_1)^2 + (Y_i - \hat{\theta}_2)^2)} = kR_i$  pro  $i = 1, \dots, n$ . A z věty 9 dostaneme, že i odhad parametru  $\beta$  (a tedy i výběrový směrový kvantil a periodický regresní kvantil) bude roven  $k$  násobku odhadu v původním výběru.

Z popsání postupu je vidět, že pro různou změnu měřítka jednotlivých složek už nemáme zaručeno, že můžeme použít větu o ekvivarianci regresního odhadu. V tomto případě již odhad není ekvivariantní.

Situace se trochu komplikuje u ekvivalence vzhledem k rotaci, kde už nevystačíme s tvrzením věty 9. Nechť náhodný výběr  $(X'_i, Y'_i)$ ,  $i = 1, \dots, n$  vznikne rotací původního výběru o úhel  $\varphi_0$ . Přejdem k polárním souřadnicím dostaneme výběr  $(R_i, F_i - \varphi_0)$ ,  $i = 1, \dots, n$ . Regresory budou v tomto případě vypadat

$$\begin{pmatrix} 1 & \cos(F_1 - \varphi_0) & \dots & \cos p(F_1 - \varphi_0) & \sin(F_1 - \varphi_0) & \dots & \sin p(F_1 - \varphi_0) \\ 1 & \cos(F_2 - \varphi_0) & \dots & \cos p(F_2 - \varphi_0) & \sin(F_2 - \varphi_0) & \dots & \sin p(F_2 - \varphi_0) \\ \vdots & & & & & & \\ 1 & \cos(F_n - \varphi_0) & \dots & \cos p(F_n - \varphi_0) & \sin(F_n - \varphi_0) & \dots & \sin p(F_n - \varphi_0) \end{pmatrix}$$

Použijme podobné značení jako v definici 8 pro vektor  $\mathbf{c} = (a_0, a_1, \dots, a_p, b_1, \dots, b_p) = (a_0, \mathbf{a}^T, \mathbf{b}^T)^T \in \mathbb{R}^{2p+1}$ . Minimalizaci (2.2) lze psát

$$\min_{\mathbf{c} \in \mathbb{R}^{2p+1}} \sum_{i=1}^n \rho_\tau(R_i - g(F_i - \varphi_0, \mathbf{c})), \quad (3.6)$$

kde

$$g(\varphi, \mathbf{c}) = a_0 + \sum_{j=1}^p a_j \cos j\varphi + b_j \sin j\varphi.$$

Použitím součtových vzorců pro funkce sin a cos dostaneme pro  $i = 1, \dots, n$

$$\begin{aligned} g(F_i - \varphi_0, \mathbf{c}) &= a_0 + \sum_{j=1}^p a_j \cos j(F_i - \varphi_0) + b_j \sin j(F_i - \varphi_0) \\ &= a_0 + \sum_{j=1}^p a_j (\cos jF_i \cos j\varphi_0 + \sin jF_i \sin j\varphi_0) \\ &\quad + b_j (\sin jF_i \cos j\varphi_0 - b_j \cos jF_i \sin j\varphi_0) \\ &= a_0 + \sum_{j=1}^p (a_j \cos j\varphi_0 - b_j \sin j\varphi_0) \cos jF_i \\ &\quad + (a_j \sin j\varphi_0 + b_j \cos j\varphi_0) \sin jF_i \\ &= a_0 + \sum_{j=1}^p u_j \cos jF_i + v_j \sin jF_i, \end{aligned}$$

kde

$$\begin{aligned} u_j &= a_j \cos j\varphi_0 - b_j \sin j\varphi_0, \quad j = 1, \dots, p \\ v_j &= a_j \sin j\varphi_0 + b_j \cos j\varphi_0, \quad j = 1, \dots, p. \end{aligned}$$

Mezi parametry  $u_j, v_j$  a  $a_j, b_j$  je vzájemně jednoznačný vztah a tak minimalizaci (3.6) lze psát

$$\min_{(a_0, \mathbf{u}^T, \mathbf{v}^T)^T \in \mathbb{R}^{2p+1}} \sum_{i=1}^n \rho_\tau(R_i - g(F_i, (a_0, \mathbf{u}^T, \mathbf{v}^T)^T)).$$

Označme  $\hat{\xi} = (\hat{a}_0, \hat{\mathbf{u}}^T, \hat{\mathbf{v}}^T)^T$  vektor, který je řešením této minimalizace. Je zřejmé, že  $g(\varphi, \hat{\xi}) = r_p(\tau|\varphi)$ , tedy je rovna směrovému kvantilu původního výběru  $(X_i, Y_i)$ ,  $i = 1, \dots, n$ . Nás ale



zajímají odhady parametrů  $a_j, b_j$ ,  $j = 1, \dots, p$  v otočeném výběru  $(X'_i, Y'_i)$ ,  $i = 1, \dots, n$ . Ty získáme ze vztahů

$$\begin{aligned}\hat{a}_j &= \hat{u}_j \cos j\varphi_0 + \hat{v}_j \sin j\varphi_0, & j = 1, \dots, p \\ \hat{b}_j &= \hat{v}_j \cos j\varphi_0 - \hat{u}_j \sin j\varphi_0, & j = 1, \dots, p.\end{aligned}$$

Položme  $\hat{\gamma} = (\hat{a}_0, \hat{a}_1, \dots, \hat{a}_p, \hat{b}_1, \dots, \hat{b}_p)^T$ , pak  $\tau$ -směrový kvantil otočeného výběru je roven  $g(\varphi, \hat{\gamma})$ . Dále platí

$$\begin{aligned}g(\varphi, \hat{\gamma}) &= \hat{a}_0 + \sum_{j=1}^p \hat{a}_j \cos j\varphi + \hat{b}_j \sin j\varphi \\ &= \hat{a}_0 + \sum_{j=1}^p (\hat{u}_j \cos j\varphi_0 + \hat{v}_j \sin j\varphi_0) \cos j\varphi + (\hat{v}_j \cos j\varphi_0 - \hat{u}_j \sin j\varphi_0) \sin j\varphi \\ &= \hat{a}_0 + \sum_{j=1}^p \hat{u}_j (\cos j\varphi_0 \cos j\varphi - \sin j\varphi_0 \sin j\varphi) + \hat{v}_j (\cos j\varphi_0 \sin j\varphi + \sin j\varphi_0 \cos j\varphi) \\ &= \hat{a}_0 + \sum_{j=1}^p \hat{u}_j \cos j(\varphi + \varphi_0) + \hat{v}_j \sin j(\varphi + \varphi_0) \\ &= g(\varphi + \varphi_0, \hat{\zeta}) = r_p(\tau | \varphi + \varphi_0).\end{aligned}$$

Tedy směrový kvantil otočeného výběru vznikne posunutím směrového kvantilu původního výběru.

Při překlopení kolem osy  $y$  je  $i$ -tý regresor tvaru

$$(1, \cos(\pi - F_i), \dots, \cos p(\pi - F_i), \sin(\pi - F_i), \dots, \sin p(\pi - F_i)).$$

To se dá přepsat do tvaru

$$(1, -\cos F_i, \dots, -\cos pF_i, \sin F_i, \dots, \sin pF_i).$$

Označíme-li  $\hat{\beta} = (\hat{a}_0, \hat{a}_1, \dots, \hat{a}_p, \hat{b}_1, \dots, \hat{b}_p)^T$  odhad vektoru  $\beta$  v původním výběru, pak odhad v překlopeném výběru bude

$$\hat{\gamma} = (\hat{a}_0, -\hat{a}_1, \dots, -\hat{a}_p, \hat{b}_1, \dots, \hat{b}_p)^T.$$

A pro odhadnutý směrový kvantil překlopeného výběru  $g(\varphi, \hat{\gamma})$  platí

$$\begin{aligned}g(\varphi, \hat{\gamma}) &= \hat{a}_0 + \sum_{j=1}^p -\hat{a}_j \cos j\varphi + \hat{b}_j \sin j\varphi \\ &= \hat{a}_0 + \sum_{j=1}^p \hat{a}_j \cos j(\pi - \varphi) + \hat{b}_j \sin j(\pi - \varphi) \\ &= g(\pi - \varphi, \hat{\beta}) = r_p(\tau | \pi - \varphi).\end{aligned}$$

Překlopení kolem přímky v obecné poloze dostaneme kombinací výše popsaných transformací.

Odvozené vztahy shrňme v následující větě.

**Věta 20.** *Periodický regresní kvantil je ekvivariantní vzhledem k posunutí, stejné změně měřítka obou složek náhodného výběru, otočení a překlopení. Tj. necht'  $\mathcal{K}_p(\tau)$  je výběrový periodický regresní  $\tau$ -kvantil výběru  $(X_i, Y_i)$ ,  $i = 1, \dots, n$ . Pak*

(i) *Pro  $\mathcal{K}'_p(\tau)$  posunutého výběru o vektor  $\mathbf{d}$  platí*

$$\mathcal{K}'_p(\tau) = \mathbf{d} + \mathcal{K}_p(\tau).$$

(ii) *Pro  $\mathcal{K}'_p(\tau)$  získané z výběru  $(kX_i, kY_i)$ ,  $i = 1, \dots, n$ , kde  $k \in \mathbb{R}$  platí*

$$\mathcal{K}'_p(\tau) = k\mathcal{K}_p(\tau).$$

(iii) *Pro  $\mathcal{K}'_p(\tau)$  otočeného výběru  $(X'_i, Y'_i)^T = A(X_i, Y_i)^T$ ,  $i = 1, \dots, n$ , kde*

$$A = \begin{pmatrix} s & -t \\ t & s \end{pmatrix}$$

*a  $s, t \in \mathbb{R}$ , platí*

$$\mathcal{K}'_p(\tau) = A\mathcal{K}_p(\tau).$$

(iv) *Pro  $\mathcal{K}'_p(\tau)$  výběru  $(X'_i, Y'_i)^T = Z(X_i, Y_i)^T$ ,  $i = 1, \dots, n$ , překlopeného kolem přímky procházející počátkem ve směru vektoru  $(u, v)^T$ , kde*

$$Z = \begin{pmatrix} v & u \\ -u & v \end{pmatrix} \begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} v & u \\ -u & v \end{pmatrix}^{-1}$$

*platí*

$$\mathcal{K}'_p(\tau) = Z\mathcal{K}_p(\tau).$$

Tvrzení věty 13 už bohužel pro výběrovou variantu obecně neplatí. Věta 8 a ekvivariance vzhledem k otočení o libovolný úhel  $\varphi_0$  nám zajišťuje, že ke „zkřížení“ nedojde pro regresory tvaru

$$\begin{aligned} \Delta(\varphi_0) &= \frac{1}{n} \sum_{i=1}^n (1, (\cos \varphi_0 - \sin \varphi_0) \cos F_i, \dots, (\cos p\varphi_0 - \sin p\varphi_0) \cos pF_i, \\ &\quad (\cos \varphi_0 + \sin \varphi_0) \sin F_i, \dots, (\cos p\varphi_0 + \sin p\varphi_0) \sin pF_i). \end{aligned}$$

Tomu ale obecně neodpovídá žádný úhel a tedy ani žádné pozorování náhodného výběru. Nicméně tyto regresory nejsou až tak vzdáleny od regresorů v (3.4) a tak ke „zkřížení“ dochází většinou jen zřídka.

### Vícerozměrný případ

Jak již bylo řečeno, budeme uvažovat trojrozměrný případ. Úvahy o konstrukci Fourierových rozvoů funkcí dvou proměnných, o tvorbě regresorů a konečně získání samotného odhadu lze přenést analogicky i do vyšších dimenzí.

Stejně jako ve dvojrozměrném případě nejprve uvedeme nějaké nezbytné základy Fourierových rozvoji funkcí dvou proměnných. Úvahy trochu více rozepíšeme, aby bylo patrné jak postupovat u výběrů vyšší dimenze než tři. Všechny citované postupy a tvrzení jsou k vidění v Kufner & Kadlec (1969).

Mějme  $Q = \{(x, y) : 0 < x < 2\pi, 0 < y < 2\pi\}$ . Uvažujme prostor  $L_2(Q)$  všech měřitelných a integrovatelných funkcí s kvadrátem na  $Q$  ( $\int_Q f^2(x, y) dx dy < \infty$ ). Skalární součin funkcí  $f, g \in L_2(Q)$  definujeme předpisem

$$\langle f, g \rangle = \int_0^{2\pi} \int_0^{2\pi} f(x, y) \overline{g(x, y)} dx dy.$$

Podobně jako u funkcí jedné proměnné potřebujeme nalézt úplný a ortonormální trigonometrický systém funkcí. Vzhledem k němu pak provedeme rozvoj funkce v řadu. V prostoru  $L_2(Q)$  je takovým systémem soustava

$$e_{mn}(x, y) = \frac{1}{2\pi} e^{i(mx+ny)}, \quad m, n \in \mathbb{Z}. \quad (3.7)$$

Pro  $e_{mn}$ ,  $m, n \in \mathbb{Z}$  tedy platí

$$\langle e_{mn}, e_{jk} \rangle = \begin{cases} 1, & \text{pokud } m = j \text{ a } n = k \\ 0, & \text{jinak.} \end{cases}$$

Fourierova řada funkce  $f \in L_2(Q)$  vzhledem k soustavě (3.7) bude mít tvar

$$f(x, y) = \frac{1}{2\pi} \sum_{m, n=-\infty}^{+\infty} c_{mn} e^{mx+ny}, \quad (3.8)$$

kde

$$c_{mn} = \langle f, e_{mn} \rangle = \frac{1}{2\pi} \int_0^{2\pi} \int_0^{2\pi} f(x, y) e^{-i(mx+ny)} dx dy.$$

Tím máme dáno vyjádření Fourierovy řady v komplexním tvaru. Kvantilová regrese nám však neumožňuje pracovat s komplexními regresory. Proto řadu v komplexním tvaru musíme převést na reálný tvar. To provedeme použitím formule

$$e^{ix} = \cos x + i \sin x$$

v (3.8). Dostáváme:

$$f(x, y) = \sum_{m, n=0}^{+\infty} \varepsilon_{mn} [\alpha_{mn} \cos mx \cos ny + \beta_{mn} \cos mx \sin ny + \gamma_{mn} \sin mx \cos ny + \delta_{mn} \sin mx \sin ny],$$

kde pro  $m, n \in \mathbb{Z}$

$$\alpha_{mn} = \frac{1}{\pi^2} \int_0^{2\pi} \int_0^{2\pi} f(x, y) \cos mx \cos ny dx dy,$$

$$\beta_{mn} = \frac{1}{\pi^2} \int_0^{2\pi} \int_0^{2\pi} f(x, y) \cos mx \sin ny dx dy,$$

$$\gamma_{mn} = \frac{1}{\pi^2} \int_0^{2\pi} \int_0^{2\pi} f(x, y) \sin mx \cos ny dx dy,$$

$$\delta_{mn} = \frac{1}{\pi^2} \int_0^{2\pi} \int_0^{2\pi} f(x, y) \sin mx \sin ny dx dy$$

a

$$\varepsilon_{mn} = \begin{cases} \frac{1}{4} & \text{pro } m = n = 0, \\ \frac{1}{2} & \text{pro } m > 0, n = 0 \text{ a pro } m = 0, n > 0, \\ 1 & \text{pro } m > 0, n > 0. \end{cases}$$

To je Fourierova řada vzhledem k soustavě funkcí

$$\cos mx \cos ny, \cos mx \sin ny, \sin mx \cos ny, \sin mx \sin ny; \quad m, n = 0, 1, 2, \dots$$

Teď již můžeme přistoupit ke konstrukci odhadu. Postup bude v mnoha ohledech podobný jako u dvojrozměrného případu. Pro trojrozměrný náhodný výběr  $(X_1^i, X_2^i, X_3^i)$ ,  $i = 1, \dots, n$  spočteme výběrový nejhlubší bod  $(\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3)$  a přejdeme k vyjádření v polárních souřadnicích, tj. k výběru  $(r_i, \varphi_1^i, \varphi_2^i)$ ,  $i = 1, \dots, n$  splňujícímu

$$\begin{aligned} X_1^i &= \hat{\theta}_1 + r_i \sin \varphi_1^i \sin \varphi_2^i, \\ X_2^i &= \hat{\theta}_2 + r_i \sin \varphi_1^i \cos \varphi_2^i, \\ X_3^i &= \hat{\theta}_3 + r_i \cos \varphi_1^i \end{aligned} \quad (3.9)$$

pro  $\varphi_1^i \in (0, \pi)$ ,  $\varphi_2^i \in [0, 2\pi)$ .

Nabízí se myšlenka odhadnout směrový kvantil  $r(\tau|\varphi_1, \varphi_2)$ , pro nějaké  $p, q \ll n$ , trigonometrickou řadou

$$\begin{aligned} g(\tau|\varphi_1, \varphi_2) &= \sum_{m=0}^p \sum_{n=0}^q (a_{mn} \cos m\varphi_1 \cos n\varphi_2 + b_{mn} \cos m\varphi_1 \sin n\varphi_2 \\ &\quad + c_{mn} \sin m\varphi_1 \cos n\varphi_2 + d_{mn} \sin m\varphi_1 \sin n\varphi_2). \end{aligned} \quad (3.10)$$

Koeficienty  $a_{mn}$ ,  $b_{mn}$ ,  $c_{mn}$ ,  $d_{mn}$  získáme pomocí kvantilové regrese pro odezvu

$$r_i, \quad i = 1, \dots, n$$

a pro regresory (za svislou čarou jsou uvedeny koeficienty odpovídající daným regresorům):

$$\begin{array}{l}
(1, \quad \cos \varphi_1^i, \quad \cos 2\varphi_1^i, \quad \dots, \quad \cos p\varphi_1^i, \quad a_{00}, a_{10}, \dots, a_{p0} \\
\quad \cos \varphi_2^i, \quad \cos 2\varphi_2^i, \quad \dots, \quad \cos q\varphi_2^i, \quad a_{01}, \dots, a_{0q} \\
\quad \sin \varphi_2^i, \quad \sin 2\varphi_2^i, \quad \dots, \quad \sin q\varphi_2^i, \quad b_{01}, \dots, b_{0q} \\
\quad \sin \varphi_1^i, \quad \sin 2\varphi_1^i, \quad \dots, \quad \sin p\varphi_1^i, \quad c_{10}, \dots, c_{p0} \\
\quad \cos \varphi_1^i \cos \varphi_2^i, \quad \cos 2\varphi_1^i \cos \varphi_2^i, \quad \dots, \quad \cos p\varphi_1^i \cos \varphi_2^i, \quad a_{11}, \dots, a_{p1} \\
\quad \cos \varphi_1^i \cos 2\varphi_2^i, \quad \cos 2\varphi_1^i \cos 2\varphi_2^i, \quad \dots, \quad \cos p\varphi_1^i \cos 2\varphi_2^i, \quad a_{12}, \dots, a_{p2} \\
\quad \vdots \\
\quad \cos \varphi_1^i \cos q\varphi_2^i, \quad \cos 2\varphi_1^i \cos q\varphi_2^i, \quad \dots, \quad \cos p\varphi_1^i \cos q\varphi_2^i, \quad a_{1q}, \dots, a_{pq} \\
\quad \cos \varphi_1^i \sin \varphi_2^i, \quad \cos 2\varphi_1^i \sin \varphi_2^i, \quad \dots, \quad \cos p\varphi_1^i \sin \varphi_2^i, \quad b_{11}, \dots, b_{p1} \\
\quad \cos \varphi_1^i \sin 2\varphi_2^i, \quad \cos 2\varphi_1^i \sin 2\varphi_2^i, \quad \dots, \quad \cos p\varphi_1^i \sin 2\varphi_2^i, \quad b_{12}, \dots, b_{p2} \\
\quad \vdots \\
\quad \cos \varphi_1^i \sin q\varphi_2^i, \quad \cos 2\varphi_1^i \sin q\varphi_2^i, \quad \dots, \quad \cos p\varphi_1^i \sin q\varphi_2^i, \quad b_{1q}, \dots, b_{pq} \\
\quad \sin \varphi_1^i \cos \varphi_2^i, \quad \sin 2\varphi_1^i \cos \varphi_2^i, \quad \dots, \quad \sin p\varphi_1^i \cos \varphi_2^i, \quad c_{11}, \dots, c_{p1} \\
\quad \sin \varphi_1^i \cos 2\varphi_2^i, \quad \sin 2\varphi_1^i \cos 2\varphi_2^i, \quad \dots, \quad \sin p\varphi_1^i \cos 2\varphi_2^i, \quad c_{12}, \dots, c_{p2} \\
\quad \vdots \\
\quad \sin \varphi_1^i \cos q\varphi_2^i, \quad \sin 2\varphi_1^i \cos q\varphi_2^i, \quad \dots, \quad \sin p\varphi_1^i \cos q\varphi_2^i, \quad c_{1q}, \dots, c_{pq} \\
\quad \sin \varphi_1^i \sin \varphi_2^i, \quad \sin 2\varphi_1^i \sin \varphi_2^i, \quad \dots, \quad \sin p\varphi_1^i \sin \varphi_2^i, \quad d_{11}, \dots, d_{p1} \\
\quad \sin \varphi_1^i \sin 2\varphi_2^i, \quad \sin 2\varphi_1^i \sin 2\varphi_2^i, \quad \dots, \quad \sin p\varphi_1^i \sin 2\varphi_2^i, \quad d_{12}, \dots, d_{p2} \\
\quad \vdots \\
\quad \sin \varphi_1^i \sin q\varphi_2^i, \quad \sin 2\varphi_1^i \sin q\varphi_2^i, \quad \dots, \quad \sin p\varphi_1^i \sin q\varphi_2^i) \quad d_{1q}, \dots, d_{pq}
\end{array} \quad (3.11)$$

pro  $i = 1, \dots, n$ .

Řada (3.10) je spojitá a  $2\pi$ -periodická v obou proměnných. Tyto vlastnosti však nestačí. Při odhadu vzhledem k regresorům (3.11) totiž opomíjíme následující dva problémy:

- (i) Jak již bylo řečeno řada (3.10) je  $2\pi$  periodická v proměnné  $\varphi_1$ . Po přechodu k polárním souřadnicím však tato úhlová veličina nabývá hodnot jen v intervalu  $(0, \pi)$ . Nabízí se tedy otázka, zda je pro hodnoty  $\varphi_1 \in (\pi, 2\pi)$  řada (3.10) definována dobře a pokud ne, tak jak ji pro tyto hodnoty definovat.
- (ii) Druhý problém se také hlavně týká této úhlové veličiny. Transformace musí být dobře definovaná, tj. každému bodu přiřadit právě jeden bod. Dále musí být prostá na otevřené množině pravděpodobnostní míry 1. Proto při transformaci vylučujeme přímku určenou osou  $x_3$ , tedy množinu  $\{(0, 0, t) : t \in \mathbb{R}\}$ . V polárních souřadnicích lze body z této množiny vyjádřit pro libovolnou hodnotu  $\varphi_2$  jako množinu  $\{(r, 0, \varphi_2) : r > 0\} \cup \{(r, \pi, \varphi_2) : r > 0\}$ . Podobně jako v předchozím bodu budeme požadovat, aby se odhad pro  $\varphi_1 = 0, \pi$  choval rozumně.

Oba dva výše popsané problémy jsou zapříčiněné vlastnostmi polárních souřadnic. Stačí si uvědomit, jak jsou data v polárních souřadnicích charakterizována. Pro bod  $(a, b, c)$ , který má ve vyjádření v polárních souřadnicích daným (3.9) souřadnice  $(r, \varphi_1, \varphi_2)$  dostaneme polohu bodu v prostoru následujícím způsobem (viz také obrázek 3.5):

Úhlová veličina  $\varphi_2$  udává o jaký úhel pootočíme kolem osy  $x_3$  polorovinu danou osami  $x_2$  a  $x_3$  s hranicí v ose  $x_3$ . V této otočené polorovině pootočíme kladnou poloosu  $x_3$  o úhel daný  $\varphi_1$ . Polohu bodu  $(a, b, c)$  najdeme na této otočené poloose ve vzdálenosti  $r$  od počátku.

Z toho pro (i) vyplývá, že bod s polárními souřadnicemi  $(r, \pi + \xi, \varphi_2)$ ,  $\xi \in (0, \pi)$ , má po přechodu ke kartézským souřadnicím stejné koordináty jako bod  $(r, \pi - \xi, \varphi_2 + \pi)$ . Rozšíříme-li tedy výběrový směrový kvantil  $r(\tau|\varphi_1, \varphi_2)$  pro hodnoty  $\varphi_1 \in (\pi, 2\pi)$  musí platit:

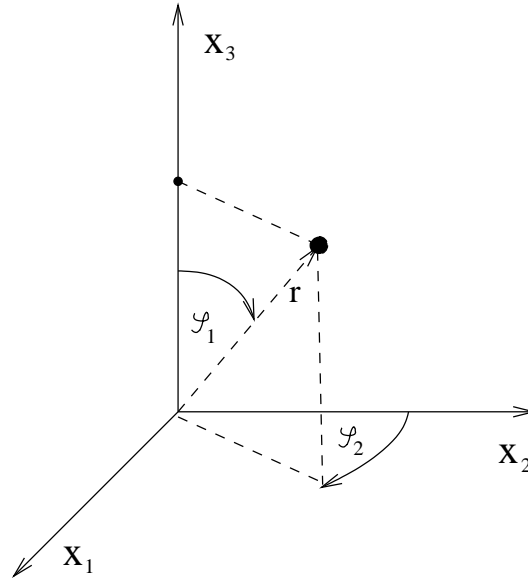
$$r(\tau|\pi + \xi, \varphi_2) = r(\tau|\pi - \xi, \varphi_2 + \pi), \quad \xi \in (0, \pi). \quad (3.12)$$

Řada (3.10) toto obecně splňovat nemusí. Podívejme se jak budou vypadat koeficienty Fourierova rozvoje funkce  $f \in L_2(Q)$ , která splňuje (3.12).

$$\begin{aligned} \alpha_{mn} &= \int_0^{2\pi} \int_0^{2\pi} f(x, y) \cos mx \cos ny \, dx \, dy \\ &= \int_0^\pi \left( \int_0^{2\pi} f(x, y) \cos mx \cos ny \, dy \right) dx \\ &\quad + \int_\pi^{2\pi} \left( \int_0^{2\pi} f(2\pi - x, y + \pi) \cos mx \cos ny \, dy \right) dx \\ &= \int_0^\pi \left( \int_0^{2\pi} f(x, y) \cos mx \cos ny \, dy \right) dx \\ &\quad - \int_\pi^0 \left( \int_0^{2\pi} f(u, y + \pi) \cos m(2\pi - u) \cos ny \, dy \right) du \\ &= \int_0^\pi \left( \int_0^{2\pi} f(x, y) \cos mx \cos ny \, dy \right) dx \\ &\quad + \int_0^\pi \left( \int_\pi^{3\pi} f(u, v) \cos mu \cos n(v - \pi) \, dv \right) du \\ &= \int_0^\pi \left( \int_0^{2\pi} f(x, y) \cos mx \cos ny \, dy \right) dx \\ &\quad + (-1)^n \int_0^\pi \left( \int_\pi^{3\pi} f(u, v) \cos mu \cos nv \, dv \right) du \\ &= (1 + (-1)^n) \int_0^\pi \int_0^{2\pi} f(x, y) \cos mx \cos ny \, dy \, dx \\ &= \begin{cases} 0, & \text{pro } n \text{ liché} \\ \neq 0, & \text{pro } n \text{ sudé.} \end{cases} \end{aligned}$$

K výpočtu jsme použili vztahy

$$\cos m(2\pi - u) = \cos mu, \quad \cos n(v - \pi) = (-1)^n \cos nv.$$



Obrázek 3.5: Reprezentace bodů v polárních souřadnicích.

Podobně dojdeme ke vztahům pro zbývající koeficienty. Platí

$$\begin{aligned}
 \alpha_{mn} &= \begin{cases} 0, & \text{pro } n \text{ liché} \\ \neq 0, & \text{pro } n \text{ sudé a rovné } 0 \end{cases} \\
 \beta_{mn} &= \begin{cases} 0, & \text{pro } n \text{ liché} \\ \neq 0, & \text{pro } n \text{ sudé a rovné } 0 \end{cases} \\
 \gamma_{mn} &= \begin{cases} 0, & \text{pro } n \text{ sudé a rovné } 0 \\ \neq 0, & \text{pro } n \text{ liché} \end{cases} \\
 \delta_{mn} &= \begin{cases} 0, & \text{pro } n \text{ sudé a rovné } 0 \\ \neq 0, & \text{pro } n \text{ liché} \end{cases}
 \end{aligned} \tag{3.13}$$

Podmínka (3.12) lze také přepsat do tvaru

$$r(\tau|\varphi_1, \varphi_2) = r(\tau|2\pi - \varphi_1, \pi + \varphi_2).$$

Úpravami řady (3.10) dostaneme

$$\begin{aligned}
 g(\tau|2\pi - \varphi_1, \pi + \varphi_2) &= \\
 &= \sum_{m=0}^p \sum_{n=0}^q (-1)^n a_{mn} \cos m\varphi_1 \cos n\varphi_2 + (-1)^n b_{mn} \cos m\varphi_1 \sin n\varphi_2 \\
 &\quad + (-1)^{n+1} c_{mn} \sin m\varphi_1 \cos n\varphi_2 + (-1)^{n+1} d_{mn} \sin m\varphi_1 \sin n\varphi_2.
 \end{aligned}$$

A je ihned vidět, že tato řada bude rovná  $g(\tau|\varphi_1, \varphi_2)$  pro koeficienty  $a_{mn}$ ,  $b_{mn}$ ,  $c_{mn}$ ,  $d_{mn}$  splňující (3.13). Pro odhad  $r(\tau|\varphi_1, \varphi_2)$  tedy použijeme jen ty regresory z (3.11), které odpovídají těmto nenulovým koeficientům.

Z popisu problému (ii) je vidět, že bod ležící na kladné poloose  $x_3$  lze v polárních souřadnicích charakterizovat nezávisle na hodnotě  $\varphi_2$  pomocí souřadnic  $(r, 0, \varphi_2)$ . Podobně bod na záporné poloose pomocí  $(r, \pi, \varphi_2)$ . Z toho pro směrový kvantil dostáváme, že  $r(\tau|0, \varphi_2)$  a  $r(\tau|\pi, \varphi_2)$  jsou konstantními funkcemi proměnné  $\varphi_2$ . Dosazením do Fourierova rozvoje  $r(\tau|\varphi_1, \varphi_2)$  pro libovolné  $\varphi_2$  dostáváme

$$\begin{aligned} r(\tau|0, \varphi_2) &= \sum_{m=0}^{+\infty} \sum_{n=0}^{+\infty} (\alpha_{mn} \cos n\varphi_2 + \beta_{mn} \sin n\varphi_2) \\ &= \sum_{n=0}^{+\infty} \left[ \left( \sum_{m=0}^{+\infty} \alpha_{mn} \right) \cos n\varphi_2 + \left( \sum_{m=0}^{+\infty} \beta_{mn} \right) \sin n\varphi_2 \right] \\ &= \sum_{m=0}^{+\infty} \alpha_{m0} + \sum_{n=1}^{+\infty} \left[ \left( \sum_{m=0}^{+\infty} \alpha_{mn} \right) \cos n\varphi_2 + \left( \sum_{m=0}^{+\infty} \beta_{mn} \right) \sin n\varphi_2 \right]. \end{aligned}$$

Zavedením parametrizačních podmínek

$$\begin{aligned} \sum_{m=0}^{+\infty} \alpha_{mn} &= 0, \quad \text{pro } n \text{ sudé a nenulové,} \\ \sum_{m=0}^{+\infty} \beta_{mn} &= 0, \quad \text{pro } n \text{ sudé a nenulové} \end{aligned}$$

dostáváme

$$r(\tau|0, \varphi_2) = \sum_{m=0}^{+\infty} \alpha_{m0}, \quad \varphi_2 \in \mathbb{R}.$$

Podobně pro  $\varphi_1 = \pi$

$$\begin{aligned} r(\tau|\pi, \varphi_2) &= \sum_{m=0}^{+\infty} \sum_{n=0}^{+\infty} (-1)^m (\alpha_{mn} \cos n\varphi_2 + \beta_{mn} \sin n\varphi_2) \\ &= \sum_{n=0}^{+\infty} \left[ \left( \sum_{m=0}^{+\infty} (-1)^m \alpha_{mn} \right) \cos n\varphi_2 + \left( \sum_{m=0}^{+\infty} (-1)^m \beta_{mn} \right) \sin n\varphi_2 \right] \\ &= \sum_{m=0}^{+\infty} (-1)^m \alpha_{m0} + \sum_{n=1}^{+\infty} \left[ \left( \sum_{m=0}^{+\infty} (-1)^m \alpha_{mn} \right) \cos n\varphi_2 + \left( \sum_{m=0}^{+\infty} (-1)^m \beta_{mn} \right) \sin n\varphi_2 \right] \end{aligned}$$

a znovu pomocí podmínek

$$\begin{aligned} \sum_{m=0}^{+\infty} (-1)^m \alpha_{mn} &= 0, \quad \text{pro } n \text{ sudé a nenulové,} \\ \sum_{m=0}^{+\infty} (-1)^m \beta_{mn} &= 0, \quad \text{pro } n \text{ sudé a nenulové} \end{aligned}$$



dostáváme konstantnost funkce  $r(\tau|\pi, \cdot)$ .

Je zřejmé, že zavedením podmínek

$$\begin{aligned} \sum_{m=0}^p a_{mn} &= 0, \\ \sum_{m=0}^p b_{mn} &= 0, \\ \sum_{m=0}^p (-1)^m a_{mn} &= 0, \\ \sum_{m=0}^p (-1)^m b_{mn} &= 0 \end{aligned} \quad (3.14)$$

dostaneme požadovanou konstantnost pro  $\varphi_1 = 0, \pi$  i pro řadu (3.10).

Zkusme se nyní podívat, jak by měly vypadat regresory pro náš odhad. U podmínek (3.13) je situace vcelku jasná. U (3.14) je postup nepatrně komplikovanější. Předpokládejme nejprve, že řád odhadu  $p$  je sudé číslo. Z podmínky (3.14) pro koeficienty  $a_{mn}$  máme

$$a_{pn} = -a_{0n} - a_{1n} - \dots - a_{p-1,n}$$

Po dosazení do druhé rovnice pro tento koeficient dostáváme

$$a_{0n} - a_{1n} + a_{2n} - a_{3n} + \dots + (-a_{0n} - a_{1n} - \dots - a_{p-1,n}) = 0$$

a po úpravě

$$a_{p-1,n} = -a_{1n} - a_{3n} - \dots - a_{p-3,n}. \quad (3.15)$$

Po dosazení do výrazu pro  $a_{pn}$  získáme

$$a_{pn} = -a_{0n} - a_{2n} - a_{4n} - \dots - a_{p-2,n}. \quad (3.16)$$

Analogické vztahy platí i pro koeficienty  $b_{mn}$ .

Sčítance z (3.10) v nichž se vyskytují koeficienty  $a_{mn}$  pro pevné nenulové  $n$  budou, po dosazení předchozích vztahů a úpravách, vypadat

$$\begin{aligned} \sum_{m=0}^p a_{mn} \cos m\varphi_1 \cos n\varphi_2 = \\ \cos n\varphi_2 \{ a_{0n}[1 - \cos p\varphi_1] + a_{1n}[\cos \varphi_1 - \cos(p-1)\varphi_1] + a_{2n}[\cos 2\varphi_1 - \cos p\varphi_1] \\ + a_{3n}[\cos 3\varphi_1 - \cos(p-1)\varphi_1] + a_{4n}[\cos 4\varphi_1 - \cos p\varphi_1] \\ + \dots + a_{p-3,n}[\cos(p-3)\varphi_1 - \cos(p-1)\varphi_1] + a_{p-2,n}[\cos(p-2)\varphi_1 - \cos p\varphi_1] \}. \end{aligned}$$

Podobně pro liché  $p$  platí

$$\begin{aligned} a_{p-1,n} &= -a_{0n} - a_{2n} - \dots - a_{p-3,n}, \\ a_{pn} &= -a_{1n} - a_{3n} - a_{5n} - \dots - a_{p-2,n} \end{aligned}$$

a

$$\begin{aligned}
& \sum_{m=0}^p a_{mn} \cos m\varphi_1 \cos n\varphi_2 = \\
& \cos n\varphi_2 \{ a_{0n}[1 - \cos(p-1)\varphi_1] + a_{1n}[\cos \varphi_1 - \cos p\varphi_1] + a_{2n}[\cos 2\varphi_1 - \cos(p-1)\varphi_1] \\
& + a_{3n}[\cos 3\varphi_1 - \cos p\varphi_1] + a_{4n}[\cos 4\varphi_1 - \cos(p-1)\varphi_1] \\
& + \dots + a_{p-3,n}[\cos(p-3)\varphi_1 - \cos(p-1)\varphi_1] + a_{p-2,n}[\cos(p-2)\varphi_1 - \cos p\varphi_1] \}. \quad (3.17)
\end{aligned}$$

Vztahy pro koeficienty  $b_{mn}$  dostaneme z předchozích záměnou  $a_{mn}$  za  $b_{mn}$  a  $\cos n\varphi_2$  za  $\sin n\varphi_2$ .

Odsud již vidíme, jak budou vypadat regresory pro náš odhad. Pro sudé  $p$  bude  $i$ -tý regresor tvaru

$$\begin{aligned}
& [ 1, \cos \varphi_1^i, \cos 2\varphi_1^i, \dots, p\varphi_1^i, \\
& (1 - \cos p\varphi_1^i) \cos 2\varphi_2^i, (\cos \varphi_1^i - \cos(p-1)\varphi_1^i) \cos 2\varphi_2^i, \dots \\
& \dots, (\cos(p-3)\varphi_1^i - \cos(p-1)\varphi_1^i) \cos 2\varphi_2^i, (\cos(p-2)\varphi_1^i - \cos p\varphi_1^i) \cos 2\varphi_2^i, \\
& (1 - \cos p\varphi_1^i) \cos 4\varphi_2^i, (\cos \varphi_1^i - \cos(p-1)\varphi_1^i) \cos 4\varphi_2^i, \dots \\
& \dots, (\cos(p-3)\varphi_1^i - \cos(p-1)\varphi_1^i) \cos 4\varphi_2^i, (\cos(p-2)\varphi_1^i - \cos p\varphi_1^i) \cos 4\varphi_2^i, \\
& \vdots \\
& (1 - \cos p\varphi_1^i) \sin 2\varphi_2^i, (\cos \varphi_1^i - \cos(p-1)\varphi_1^i) \sin 2\varphi_2^i, \dots \\
& \dots, (\cos(p-3)\varphi_1^i - \cos(p-1)\varphi_1^i) \sin 2\varphi_2^i, (\cos(p-2)\varphi_1^i - \cos p\varphi_1^i) \sin 2\varphi_2^i, \\
& (1 - \cos p\varphi_1^i) \sin 4\varphi_2^i, (\cos \varphi_1^i - \cos(p-1)\varphi_1^i) \sin 4\varphi_2^i, \dots \\
& \dots, (\cos(p-3)\varphi_1^i - \cos(p-1)\varphi_1^i) \sin 4\varphi_2^i, (\cos(p-2)\varphi_1^i - \cos p\varphi_1^i) \sin 4\varphi_2^i, \quad (3.18) \\
& \vdots \\
& \sin \varphi_1^i \cos \varphi_2^i, \sin 2\varphi_1^i \cos \varphi_2^i, \dots, \sin p\varphi_1^i \cos \varphi_2^i, \\
& \sin \varphi_1^i \cos 3\varphi_2^i, \sin 2\varphi_1^i \cos 3\varphi_2^i, \dots, \sin p\varphi_1^i \cos 3\varphi_2^i, \\
& \vdots \\
& \sin \varphi_1^i \sin \varphi_2^i, \sin 2\varphi_1^i \sin \varphi_2^i, \dots, \sin p\varphi_1^i \sin \varphi_2^i, \\
& \sin \varphi_1^i \sin 3\varphi_2^i, \sin 2\varphi_1^i \sin 3\varphi_2^i, \dots, \sin p\varphi_1^i \sin 3\varphi_2^i, \\
& \vdots \\
& ]
\end{aligned}$$

To odpovídá koeficientům

$$\begin{aligned}
 & [ a_{00}, a_{10}, a_{20}, \dots, a_{p0}, \\
 & a_{02}, a_{12}, \dots, a_{p-3,2}, a_{p-2,2}, \\
 & a_{04}, a_{14}, \dots, a_{p-3,4}, a_{p-2,4}, \\
 & \vdots \\
 & b_{02}, b_{12}, \dots, b_{p-3,2}, b_{p-2,2}, \\
 & b_{04}, b_{14}, \dots, b_{p-3,4}, b_{p-2,4}, \\
 & \vdots \\
 & c_{11}, c_{21}, \dots, c_{p1}, \\
 & c_{13}, c_{23}, \dots, c_{p3}, \\
 & \vdots \\
 & d_{11}, d_{21}, \dots, d_{p1}, \\
 & d_{13}, d_{23}, \dots, d_{p3}, \\
 & \vdots \qquad \qquad \qquad ] \tag{3.19}
 \end{aligned}$$

Regresory pro liché  $p$  snadno získáme z výrazu (3.17) a tak je již nebudeme vypisovat.

**Poznámka.** Všimněme si, že soustava (3.18) je ortogonální, ale již není úplná. Označíme-li  $h_{ij}$  členy této soustavy, pak Fourierovy koeficienty směrového kvantilu  $r$  budou zřejmě rovny  $\langle r, h_{ij} \rangle$ .

Označme vektor v (3.18)  $\mathbf{h}$  a dívejme se na něj jako funkci dvou proměnných  $\varphi_1, \varphi_2$ . Dále označme, pro  $\tau \in (0, 1)$ , odhad vektoru (3.19) získaný kvantilovou regresí pro regresory dané v (3.18) a odezvu  $r_i, i = 1, \dots, n$  jako  $\hat{\beta}(\tau)$ . Pak analogicky jako v dvojrozměrném případě definujme trojrozměrný  $\tau$ -směrový kvantil jako

$$r_{pq}(\tau|\varphi_1, \varphi_2) = \sum_i \hat{\beta}_i(\tau) h_i(\varphi_1, \varphi_2).$$

Dále pak trojrozměrný výběrový periodický regresní  $\tau$  kvantil  $\mathcal{K}_{pq}(\tau)$  jako uzávěr množiny

$$\begin{aligned}
 \{(x_1, x_2, x_3) : \quad x_1 &= \hat{\theta}_1 + r_{pq}(\tau|\varphi_1, \varphi_2) \sin \varphi_1 \sin \varphi_2, \\
 x_2 &= \hat{\theta}_2 + r_{pq}(\tau|\varphi_1, \varphi_2) \sin \varphi_1 \cos \varphi_2, \\
 x_3 &= \hat{\theta}_3 + r_{pq}(\tau|\varphi_1, \varphi_2) \cos \varphi_1, \\
 &\varphi_1 \in [0, \pi], \varphi_2 \in [0, 2\pi)\}
 \end{aligned}$$

Věty 18, 19 platí i pro trojrozměrný kvantil. Ekvivariance vzhledem k posunutí a stejné změně měřítka se ukáže analogicky jako ve dvojrozměrném případě (věta 20).

Ekvivariance vzhledem k otočení platí jen při rotaci kolem osy  $x_3$ , tedy rotaci danou změnou úhlu  $\varphi_2$ . Ukáže se to podobným způsobem jako v dvojrozměrném případě. Při otočení

náhodného výběru kolem osy  $x_3$  o úhel  $\xi$  přejdeme k náhodnému výběru, který má polární souřadnice  $(r_i, \varphi_1^i, \varphi_2^i - \xi)$ ,  $i = 1, \dots, n$ . Vezmeme  $m = 1$ ,  $n = 2$ , pak  $c_{12} = d_{12} = 0$ . Zaměříme se jen na sčítance pro tyto hodnoty  $m$ ,  $n$ . Dostáváme

$$\begin{aligned} & a_{12} [\cos \varphi_1^i - \cos(p-1)\varphi_1^i] \cos 2(\varphi_2^i - \xi) + b_{12} [\cos \varphi_1^i - \cos(p-1)\varphi_1^i] \sin 2(\varphi_2^i - \xi) \\ &= (a_{12} \cos 2\xi - b_{12} \sin 2\xi) [\cos \varphi_1^i - \cos(p-1)\varphi_1^i] \cos 2\varphi_2^i \\ & \quad + (a_{12} \sin 2\xi + b_{12} \cos 2\xi) [\cos \varphi_1^i - \cos(p-1)\varphi_1^i] \sin 2\varphi_2^i \\ &= a'_{12} [\cos \varphi_1^i - \cos(p-1)\varphi_1^i] \cos 2\varphi_2^i + b'_{12} [\cos \varphi_1^i - \cos(p-1)\varphi_1^i] \sin 2\varphi_2^i. \end{aligned}$$

Z tohoto by, analogicky jako v dvojrozměrném případě, měla být dokazovaná vlastnost patrná. Pro ostatní hodnoty  $m, n$  je postup stejný.

Ještě bychom se ale měli přesvědčit, zda koeficienty  $a'_{mn}, b'_{mn}, c'_{mn}, d'_{mn}$  splňují podmínky (3.13) a (3.14). O tom se přesvědčíme snadno. Nulovost  $a_{mn}, b_{mn}$  implikuje nulovost  $a'_{mn}, b'_{mn}$ . Podobně tak u ostatních koeficientů. Dále

$$\sum_{m=0}^p a'_{mn} = \cos n\xi \sum_{m=0}^p a_{mn} - \sin n\xi \sum_{m=0}^p b_{mn} = 0.$$

Znovu si vystačíme s tvrzením, že u ostatních koeficientů a druhé podmínky v (3.14) budeme postupovat zcela analogicky.

Vzhledem k tomu, že  $r_{pq}(\tau|0, \cdot)$  a  $r_{pq}(\tau|\pi, \cdot)$  jsou konstantními funkcemi, tak ihned můžeme vyloučit, že by odhad byl ekvivariantní vůči rotaci v úhlové veličině  $\varphi_1$ .

Ještě předtím než uvedeme konkrétní příklady výběrových periodických kvantilů pro určitá rozdělení se podíváme na metody, které mohou nějakým způsobem zlepšit jejich odhad.

### 3.3 Metody pro zlepšení odhadu výběrových periodických kvantilů

#### 1. Normalizace měřítka

Mějme dva body v rovině, pro které po transformaci do polárních souřadnic platí, že vzdálenost obou bodů od počátku je rovna nějakému  $r$ . Rozdíl úhlů určující tyto body označme  $\alpha$ . Je zřejmé, že s rostoucím  $r$  se bude vzdálenost těchto dvou bodů zvětšovat, přičemž v grafickém vyjádření, kde souřadnicemi těchto bodů jsou vzdálenost od počátku a úhel, bude jejich vzdálenost stejná. Podobně dva body blízko středu svírají větší úhel než stejně vzdálené dva body dále od středu. Tato vlastnost se poměrně nepříjemně promítne při transformaci odhadu směrového kvantilu  $r_p(\tau|\varphi)$  zpět do kartézských souřadnic. Nepatrné zvlnění v nízkých hodnotách  $r_p$  se po transformaci projeví jako poměrně výrazné. Naopak pro vysoké hodnoty  $r_p$  bude i více patrné zvlnění po transformaci méně výrazné. Poznamenejme ještě, že není důležitá velikost hodnot  $r_p$ , ale rozdíly v jejich hodnotách. Proto bude výběrový periodický regresní kvantil pro hodnoty blíže středu mnohem výrazněji zvlněný než pro hodnoty dále od středu. Nabízí se tedy hledat takové transformace pozorovaných dat, pro které by byl směrový kvantil podobný konstantní funkci. Najít takovou transformaci je ve většině případů téměř nemožné. Proto je výhodná alespoň následující jednoduchá transformace.

*Pokud je výběrová varianční matice náhodného výběru různá od jednotkové matice, může být*

výhodné přejít vhodnou lineární transformací k výběru, jehož výběrová varianční matice je jednotková. Pokud taková lineární transformace neexistuje, přejdeme alespoň k výběru jehož složky budou mít jednotkový rozptyl. Tedy k výběru, kde pozorování jeho  $k$ -té složky, pro  $k = 1, \dots, p$ , mají tvar:

$$\tilde{X}_k^i = X_k^i / \hat{\sigma}_k, \quad i = 1, \dots, n,$$

$\hat{\sigma}_k^2$  značí nějaký odhad rozptylu  $k$ -té složky náhodného vektoru.

To je výhodné zejména u elipticky souměrných rozdělání (např. normální rozdělání). Směrový kvantil, který má po transformaci do kartézských souřadnic eliptický tvar, má totiž Fourierův rozvoj s nekonečným počtem sčítanců. Po normalizaci měřítka má periodický regresní kvantil tvar kružnice a to odpovídá konstantnímu směrovému kvantilu a jeho odhadu jen absolutním členem.

## 2. Lineární transformace

Další možností je použít na data lineární transformace vůči kterým nejsou výběrové periodické regresní kvantily ekvivariantní. Jedná se hlavně o přechod k jiné než ortonormální soustavě souřadnic. Tím dojde v některých místech k roztažení rozvržení dat a v jiných k zúžení rozvržení dat. Většinou platí následující pravidlo, jehož platnost je způsobena změnou poměru vzdáleností od středu v těchto oblastech.

*V místech, kde vlivem transformace dat dochází k roztažení rozvržení dat je nový odhad oproti původnímu více zvlněný. A naopak v místech, kde dochází k zúžení rozvržení dat, bude odhad méně zvlněný.*

## 3. Regresory pro symetrický periodický regresní kvantil

Uvažujme dvojrozměrný případ a periodický regresní kvantil osově souměrný kolem přímky procházející nejhlubším bodem. Díky vlastnosti ekvivariance vzhledem k otočení, můžeme výběr otočit tak, aby body ležící na této přímce měly po přechodu k polárním souřadnicím veličinu udávající úhel rovnou nule. Směrový kvantil  $r(\tau|\varphi)$  bude v tomto případě sudou funkcí. Pro sinové koeficienty Fourierova rozvoje platí

$$\beta_j = \frac{1}{\pi} \int_0^{2\pi} r(\tau|\varphi) \sin(j\varphi) d\varphi = 0, \quad j = 1, 2, \dots$$

Dá se tedy předpokládat, že by i odhadnuté sinové koeficienty byly blízké nule a proto směrový kvantil budeme odhadovat trigonometrickou řadou

$$r_p(\tau|\varphi) = \hat{a}_0 + \sum_{j=1}^p \hat{a}_j \cos j\varphi.$$

Při menších rozsazích výběrů je to výhodné, protože nám to umožní menší navýšení řádu  $p$ . Pro velké rozsahy je to také výhodné, neboť pro získání stejně kvalitního odhadu nám stačí poloviční počet regresorů. A výpočet je tedy méně časově i paměťově náročný.

Před použitím tohoto postupu bychom se měli přesvědčit o symetrii náhodného výběru. Pro náhodný výběr  $(X_1^i, X_2^i)$ ,  $i = 1, \dots, n$  z rozdělání daného distribuční funkcí  $F$  můžeme, po vhodném přetočení a posunutí, použít nějaký z neparametrických testů symetrie kolem přímky  $x_2 = x_1$ , tedy test hypotézy  $F(x_1, x_2) = F(x_2, x_1)$ . Např. Wilcoxonův nebo znaménkový test symetrie náhodného výběru  $Z_i = X_1^i - X_2^i$  kolem nuly.

Ve více rozměrech je situace podobná. Uvedme ještě jak bychom postupovali ve trojrozměrném případě. Po vhodném přetočení dojdeme k vektoru s rozdělením, jehož periodický regresní kvantil je zrcadlově symetrický kolem roviny určené osou  $x_3$  a nejhlubším bodem. Pak již jen vektor otočíme kolem  $x_3$  tak, aby body ležící v této rovině měly nulovou úhlovou veličinu  $\varphi_2$ . Nyní je směrový kvantil pro pevné  $\varphi_1$  sudou funkcí v proměnné  $\varphi_2$ . Pro koeficienty Fourierovy rozvoje, v nichž se vyskytují členy  $\sin j\varphi_2$ , platí

$$\begin{aligned}\beta_{mn} &= \frac{1}{\pi^2} \int_0^{2\pi} \int_0^{2\pi} r(\tau|\varphi_1, \varphi_2) \cos m\varphi_1 \sin n\varphi_2 \, d\varphi_1 \, d\varphi_2 \\ &= \frac{1}{\pi^2} \int_0^{2\pi} \left( \int_0^{2\pi} r(\tau|\varphi_1, \varphi_2) \sin n\varphi_2 \, d\varphi_2 \right) \cos m\varphi_1 \, d\varphi_1 \\ &= \frac{1}{\pi^2} \int_0^{2\pi} 0 \cos m\varphi_1 \, d\varphi_1 = 0, \\ \delta_{mn} &= \frac{1}{\pi^2} \int_0^{2\pi} \int_0^{2\pi} r(\tau|\varphi_1, \varphi_2) \sin m\varphi_1 \sin n\varphi_2 \, d\varphi_1 \, d\varphi_2 = 0.\end{aligned}$$

Tedy podobně jako v dvojrozměrném případě budeme hledat jen odhady koeficientů  $\alpha_{mn}$  a  $\gamma_{mn}$ .

### 3.4 Volba řádu $p$ a konzistence výběrového periodického regresního kvantilu

O volbě řádu již bylo něco řečeno v části pojednávající o dvojrozměrných periodických regresních kvantilech. Požadavek (3.5) na rychlost konvergence řádu  $p$  k nekonečnu lze použít také u vícerozměrných kvantilů. A stejně jako u dvojrozměrného případu nám zaručí, že pro počet bodů  $N$  ležících uvnitř výběrového kvantilu platí

$$\frac{N}{n} \xrightarrow{n \rightarrow \infty} \tau.$$

Je jasné, že velikost řádu  $p$  (resp.  $p, q$ ) budeme volit nejen v závislosti na počtu pozorování, ale také na předpokládaném rozdělení z kterého výběr pochází a dokonce také na hodnotě parametru  $\tau$ . Většinou platí, že s rostoucím  $\tau$  volíme i větší velikost řádu  $p$  (viz např. obr. 3.9). U rozdělení, která mají jen po částech hladký kvantil (Exponenciální rozdělení) bude potřeba zvýšit řád odhadu, abychom lépe vystihli tvar kvantilu v těchto bodech. To ovšem bude na úkor celého odhadu, který bude v ostatních bodech více zvlněný. Proto se dá těžko určit nějaké pravidlo pro volbu řádu odhadu. Uvedeme tabulku s doporučeními pro volbu řádu pro různé rozsahy výběru. Pro zjednodušení ve trojrozměrném případě klademe  $p = q$ .

Rozsah výběru	500	1000	5000	10000	50000
Dvojrozměrný případ	2–9	3–15	6–20	7–24	10–30
Trojrozměrný případ	1–4	2–6	3–10	4–13	6–18

Tabulka 3.1: Volba řádu odhadu  $p$  pro různé rozsahy výběru.

Ještě se ve zkratce zmiňme o možnosti důkazu konzistence. Vyjdeme z modelu (2.3) v kvantilové regresi, který vlastně používáme k našemu odhadu. Předpokládejme, že máme náhodný

výběr splňující pro  $\tau \in (0, 1)$  po přechodu do polárních souřadnic

$$r_i = r(\tau|\varphi_i) + e_i, \quad i = 1, \dots, n,$$

kde  $e_i$ ,  $i = 1, \dots, n$  jsou nezávislé, ale nemusí být stejně rozdělené. Použijeme-li značení z části 3.1, pak  $e_i$  má rozdělení dané distribuční funkcí

$$F_i(x) = Q(x + r(\tau|\varphi_i) | \varphi_i).$$

Tedy platí podmínka  $F_i(0) = \tau$  z modelu (2.3).

Uvažujeme nějaký trigonometrický systém a rozvíňme  $r(\tau|\cdot)$  ve Fourierovu řadu vzhledem k tomuto systému. Nekonečný vektor koeficientů tohoto rozvoje označme  $\beta(\tau)$ . Vektor získaný odhadem v modelu

$$r_i = \mathbf{b}^T \mathbf{x}_i^{(l)} + e_i, \quad i = 1, \dots, n,$$

označme  $\hat{\beta}_l(\tau)$ .  $l$  značí řád odhadu a  $\mathbf{x}_i^{(l)}$  značí část trigonometrického systému pro hodnotu  $\varphi_i$ , vzhledem ke kterému jsme prováděli Fourierův rozvoj, mající  $l$  prvků (např. po vhodném přeuspořádání prvních  $l$  prvků). Jde-li  $n$  k nekonečnu, půjde i řád odhadu  $l$  k nekonečnu.

Princip důkazu konzistence v Bantli & Hallin (1999) a Oberhofer (1982) pak zřejmě půjde, s nemalými technickými potížemi, přenést i na zde popsany případ. Dostali bychom

$$\hat{\beta}_l(\tau) \xrightarrow[n \rightarrow \infty]{P} \beta(\tau),$$

tedy konvergenci v pravděpodobnosti  $r_l(\tau|\varphi)$  k  $r(\tau|\varphi)$ .

## 3.5 Pár poznámek

### 1. Proč nejhlubší bod?

Hloubka bodu v jistém smyslu vypovídá o počtu pozorování v nejméně příznivém směru od tohoto bodu (poloprostor v jehož hranici leží tento bod a v němž leží nejmenší počet pozorování). Nejhlubší bod je takový, pro který v nejméně příznivém směru leží nejvíce pozorování. Toto je velice výhodné při odhadu směrového kvantilu, který dostáváme z dat transformovaných do polárních souřadnic. Vezmeme-li totiž dělení intervalu  $[0, 2\pi]$  takové, že jednotlivé úseky dělení jsou stejně dlouhé, pak vzhledem k popsané vlastnosti nejhlubšího bodu, bude počet bodů v jednotlivých úsecích (body, jejichž úhlová souřadnice padne do tohoto úseku) pro většinu „rozumných“ rozdělení víceméně podobný. Tedy, že i v tom úseku kam padne nejméně pozorování, jsme na tom stále lépe, než kdybychom zvolili jiný než nejhlubší bod.

Na druhou stranu, díky vlastnosti polárních souřadnic popsané v úvodu části 3.3, může být rozložení bodů v polárních souřadnicích pro nějaké hodnoty úhlu  $\varphi$  v jistém smyslu „řídké“ a pro jiné „zahuštěné“ (např. exponenciální rozdělení, úsek kolem  $\varphi = \frac{5}{4}\pi$  - obr. 3.7). A to svádí k myšlence, že odhad  $r_p(\tau|\varphi)$  může být pro nějaké hodnoty  $\varphi$  méně přesný.

Ale toto řídké rozložení dat v polárních souřadnicích pro určité hodnoty úhlu  $\varphi$  je způsobeno spíše větším rozptylem vzdálenosti od počátku  $r$  pro tyto hodnoty úhlu. Tedy pozorování jsou více rozházena ve směru vzdálenosti  $r$  a ne v úhlové veličině  $\varphi$ . A vzhledem k tvrzení věty 10 můžeme tedy očekávat, že bude odhad v oblastech s „řídkým“ rozložením stejně přesný jako v oblastech s „hustým“ rozložením dat.

**Poznámka.** (Úhlový nejhlubší bod). *Jistě zajímavou variantou volby středu pro transformaci do polárních souřadnic by mohla být modifikace nejhlubšího bodu, nazvěme ji třeba úhlový nejhlubší bod. Ve dvojrozměrném případě definujme úhlovou hloubku bodu, pro nějaké pevně zvolené  $\varphi_0$ , následujícím způsobem:*

*Představme si výseč určenou dvěma polopřímkami, které svírají nějaký pevně zvolený úhel  $\varphi_0$  a mají počátek v tomto bodě. Úhlovou hloubkou nazvěme počet bodů ležících v takové výseči, která obsahuje nejmenší počet bodů. Úhlovým nejhlubším bodem nazvěme bod s největší úhlovou hloubkou. Úhel  $\varphi_0$  zvolíme v závislosti na rozsahu výběru a typu rozdělení.*

*Tento postup se dá zobecnit i na prostory vyšších dimenzí. Např. v trojrozměrném případě bychom místo výseče použili rotační kužel.*

*Takto definovaný nejhlubší bod nám zřejmě zaručí ještě rovnoměrnější rozvržení bodů po přechodu do polárních souřadnic (ve smyslu bodu 1) a tedy i lepší odhad.*

## 2. Robustnost

Velkou výhodou periodických regresních kvantilů je robustnost. Při jejich konstrukci využíváme jen robustních statistických odhadů. O robustnosti nejhlubšího bodu vypovídá jeho bod zlomu - viz věta 5. Např. ve dvojrozměrném případě můžeme nejméně třetinu pozorování nahradit nekonečnou hodnotou než výběrový nejhlubší bod „ulétne“ do nekonečna. To značí velmi malou citlivost na odlehlá pozorování.

Po přechodu k polárním souřadnicím, kde za střed bereme právě nejhlubší bod, použijeme k našemu odhadu další velmi robustní metodu - kvantilovou regresi. O ní víme (viz kapitola 2), že její influenční funkce je omezená ve veličině odpovídající odezvě. Té v polárních souřadnicích odpovídají veličiny  $r_i$  vzdálenosti od počátku. Odlehlost pozorování se po přechodu k polárním souřadnicím projeví jen právě na této veličině. Fakt, že kvantilová regrese může být citlivá na odlehlé hodnoty regresorů (neomezená influenční funkce v této veličině) není pro nás významný, jelikož regresory používané k našemu odhadu vznikají z omezených funkcí  $\cos$  a  $\sin$ . Navíc nám věta 10 zaručuje, že výběrový periodický regresní kvantil se nezmění budeme-li libovolný počet pozorování, která leží mimo něj posunovat o libovolnou vzdálenost dále ve směru od nejhlubšího bodu.

Dohromady tedy získáváme velice robustní odhad, který je málo citlivý i dokonce na větší počet odlehlých pozorování.

## 3. Výpočetní a paměťová náročnost

O výpočetní časové náročnosti nejhlubšího bodu již bylo pojednáno v sekci 1.1. Paměťové nároky algoritmů DEEPLOC a HALFMED, které používáme, nepotřebují o moc více paměti než kolik zabírají jejich vstupní data. Pro menší rozsahy dat (do 5000) oba algoritmy, na počítači s 2 GHz procesorem a 512 Mb operační paměti, napočítaly nejhlubší bod téměř okamžitě. Pro trojrozměrný výběr o rozsahu 80000, výpočet algoritmem DEEPLOC už ale trval přibližně 5 hodin.

Kvantilová regrese většinou používá simplexového algoritmu. Výpočet je tedy poměrně rychlý. Konstrukce výběrových kvantilů však vyžaduje odhad velkého počtu parametrů (např. pro trojrozměrná data až několik stovek) a tedy výpočet může být poměrně časově náročný. V programu R jsou proto k dispozici další algoritmy, které provádí výpočet rychleji (ale někdy jsou jen přibližné).



## 3.6 Příklady výběrových kvantilů

V této části ukážeme jak mohou vypadat výběrové periodické regresní kvantily pro výběry vygenerované z exponenciálního a normálního rozdělení. Nakonec uvedeme jeden menší příklad aplikace na reálná data.

### Exponenciální rozdělení – dvojrozměrný případ

V následujících příkladech byl generován náhodný výběr  $(x_i, y_i)$ , kde  $x_i, y_i$  jsou nezávislé a mají exponenciální rozdělení s parametrem 1.

Obr. 3.6 – 3.10 jsou zhotoveny pro rozsah výběru 500. Pro dvojrozměrná data můžeme takový rozsah považovat za velmi malý. I přesto vidíme, že dostáváme poměrně pěkný odhad, který si poradí i s oblastmi, kde není teoretický kvantil hladký.

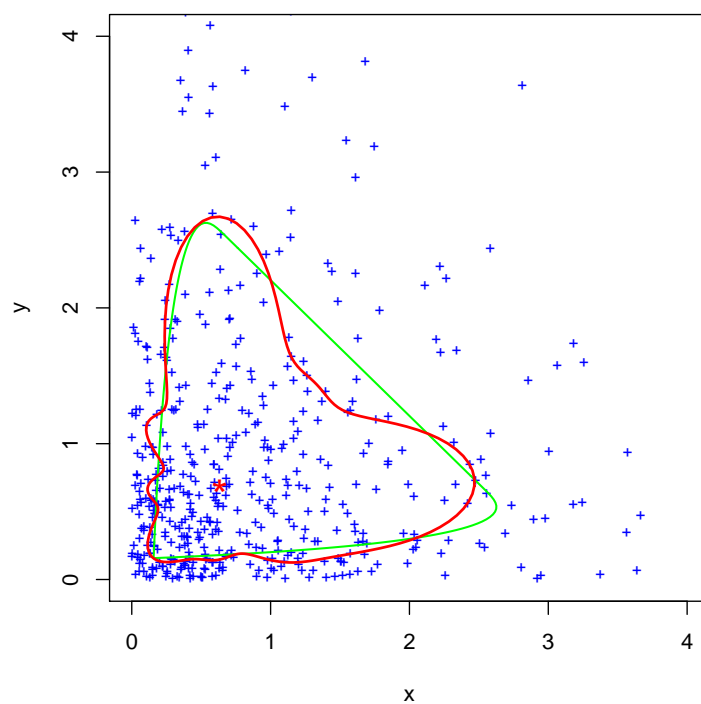
Na většině obrázcích je vidět, že kvantil je více zvlněný v oblastech, kde je blíže středu (nejhlubšímu bodu). To je způsobeno mechanismy, které jsou popsány v bodu 1 části 3.3.

Obr. 3.8 ukazuje dvě krajní, ale ještě přijatelné volby řádu odhadu. Kvantil na levém obrázku je příliš zakulacený, na pravém příliš zvlněný.

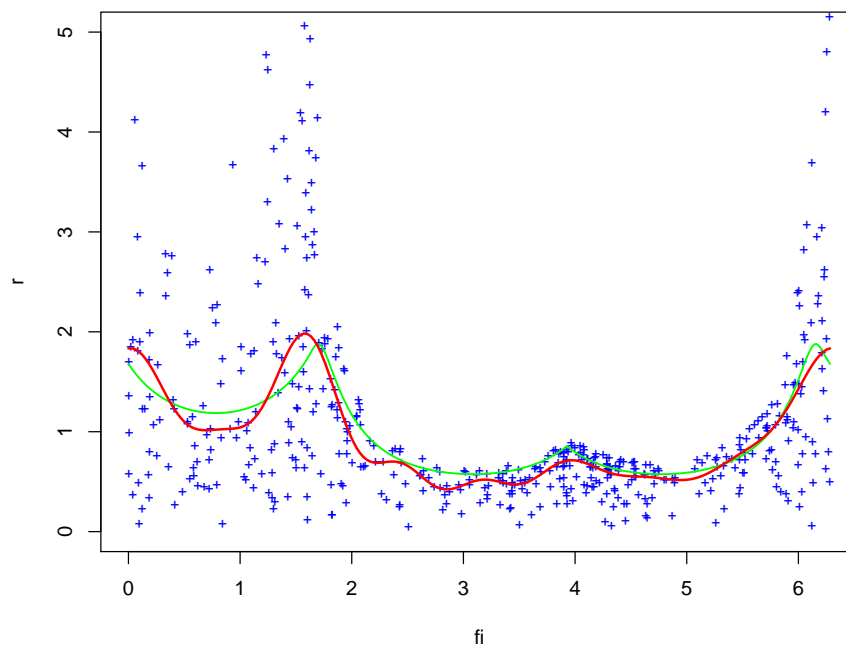
Na obr. 3.9 je zase vidět nutnost volby vyššího řádu odhadu s rostoucím  $\tau$ .

Pro kvantily na obr. 3.10 byl navíc použit odhad pro symetrický kvantil. Ten je symetrický kolem přímky  $y = x$ . Směrový kvantil je tedy symetrický kolem úhlu  $\frac{5}{4}\pi$ . Po rotaci dat o tento úhel se stane sudou funkcí. A platí pro něj vlastnosti popsané v bodě 3 části 3.3. Získaný odhad je vyobrazen na levém obrázku. Oproti odhadu bez symetrie došlo k příjemnému zlepšení. Pravý obrázek ukazuje vliv lineární transformace - data byla nejprve transformována maticí se sloupci  $(5, 2)^T$  a  $(2, 5)^T$  (to jsou zároveň směry na které přechází zobrazení os  $x$  a  $y$ ). Potom, stejně jako předtím, byl použit odhad pro symetrický kvantil. Vidíme, že došlo k jevům popsaných v bodě 2 části 3.3. Ještě dodejme, že byl odhadován 95% kvantil, kterému pro takto malý rozsah odpovídá přibližně 475 pozorování ležících uvnitř.

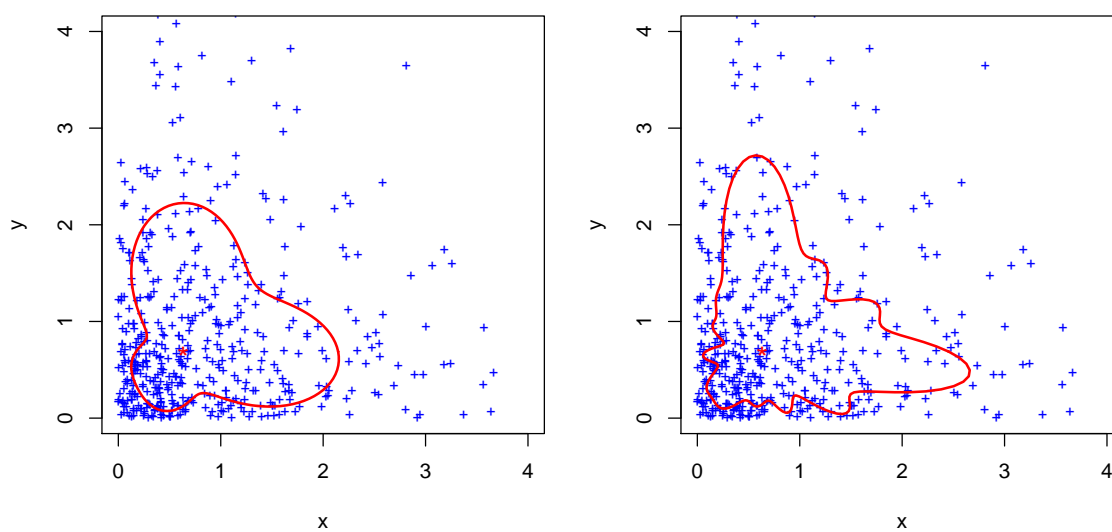
Na obr. 3.11 jsou k vidění odhadnuté kvantily pro rozsah výběru 50000.



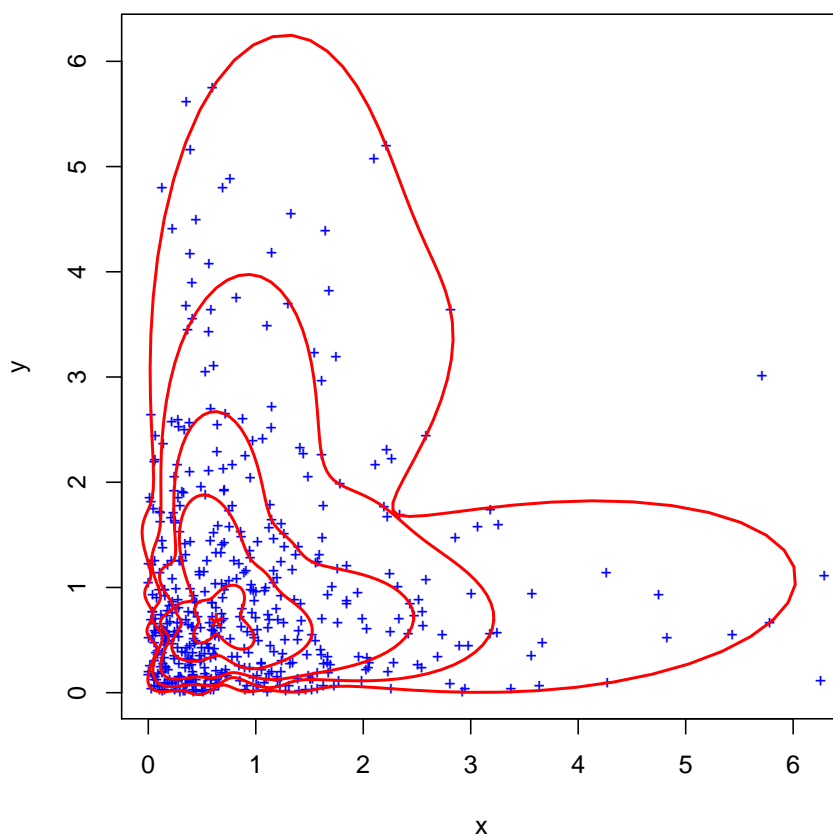
Obrázek 3.6: Výběrový a teoretický 0.5–kvantil pro exponenciální rozdělení. Rozsah výběru 500, řád rozvoje 9.



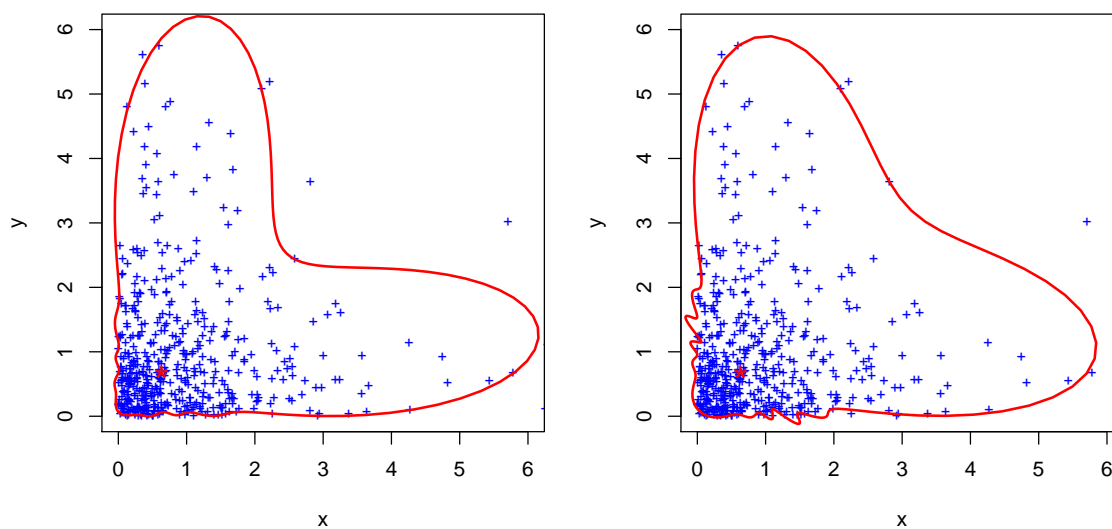
Obrázek 3.7: Výběrový a teoretický směrový 0.5–kvantil pro exponenciální rozdělení v polárních souřadnicích. Rozsah výběru 500, řád rozvoje 9.



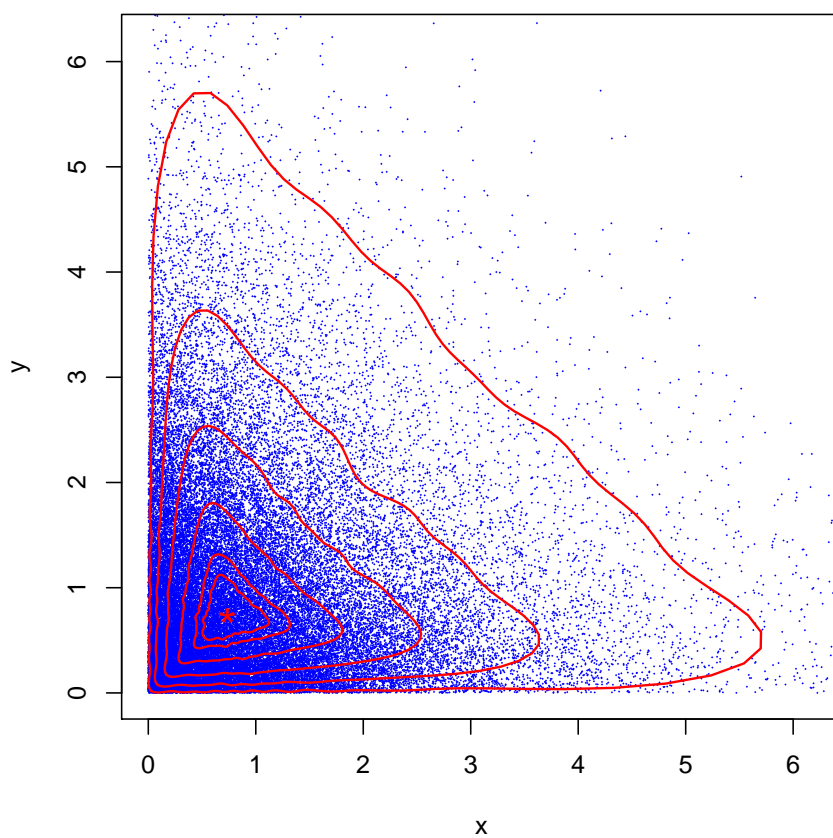
Obrázek 3.8: Výběrové 0.5-quantily pro řád 4 a řád 13. Rozsah výběru 500.



Obrázek 3.9: Výběrové kvantily pro  $\tau = 0.05, 0.25, 0.5, 0.75$  a  $0.95$ . Řád volen 4, 7, 9, 9 a 11. Exponenciální rozdělení. Rozsah výběru 500.

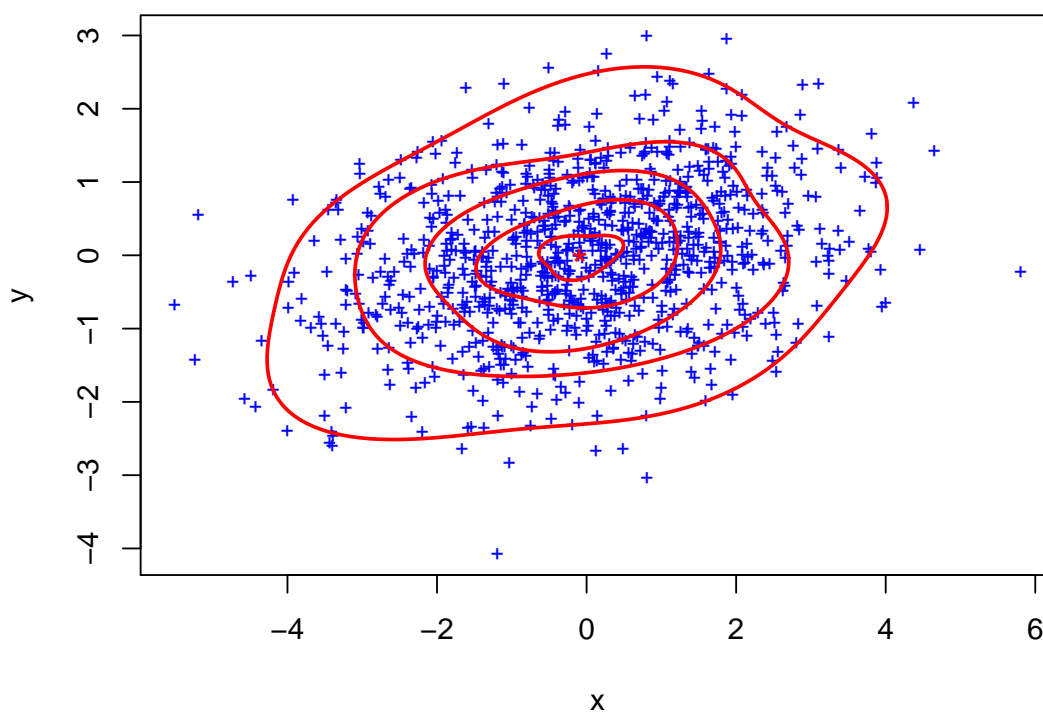


Obrázek 3.10: Výběrové 0.95-kvantily. Použity odhady pro symetrický kvantil. Napravo navíc data před odhadem prošla lineární transformací. Řád volen 11. Rozsah výběru 500.



Obrázek 3.11: Výběrové kvantily pro  $\tau = 0.05, 0.1, 0.25, 0.5, 0.75$  a  $0.95$ . Řád volen 17, 20, 24, 22, 24 a 25. Exponenciální rozdělení. Rozsah výběru 50000.

## Normální rozdělení – dvojrozměrný případ



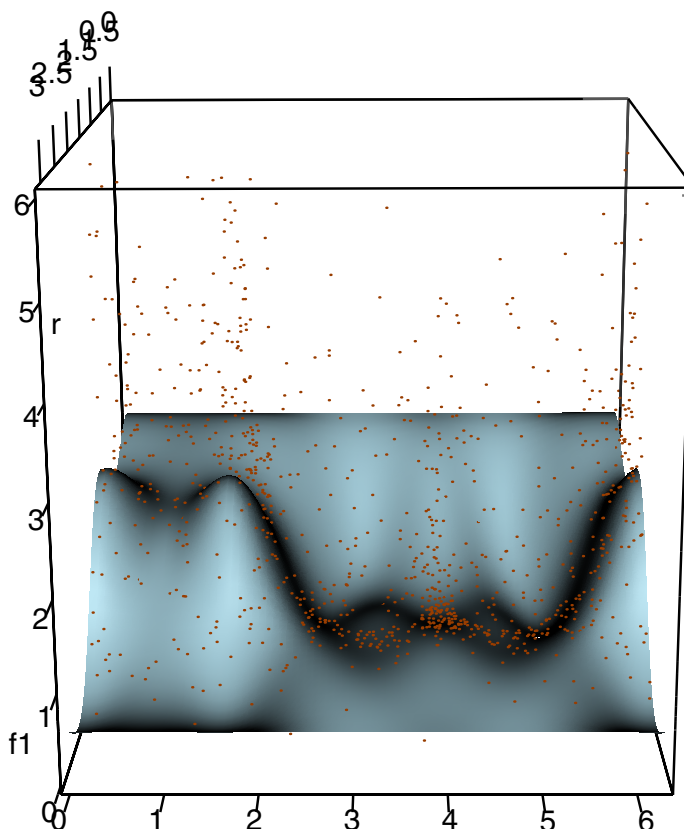
Obrázek 3.12: Výběrové kvantily pro výběr z dvojrozměrného normálního rozdělení s nulovou střední hodnotou a variační maticí s 3 a 1 na diagonále a mimodiagonálními prvky 0.5. Velikost výběru 1000. Pro odhad byl použit postup bodu 1 z sekce 3.3 (přešli jsme k výběru s jednotkovou výběrovou varianční maticí pomocí Choleského rozkladu výběrové matice původního výběru). Pro hodnoty  $\tau$  rovné 0.05, 0.1, 0.25, 0.5, 0.75 a 0.95 volen řád 3, 3, 4, 5 a 5.

### Exponenciální rozdělení – trojrozměrný případ

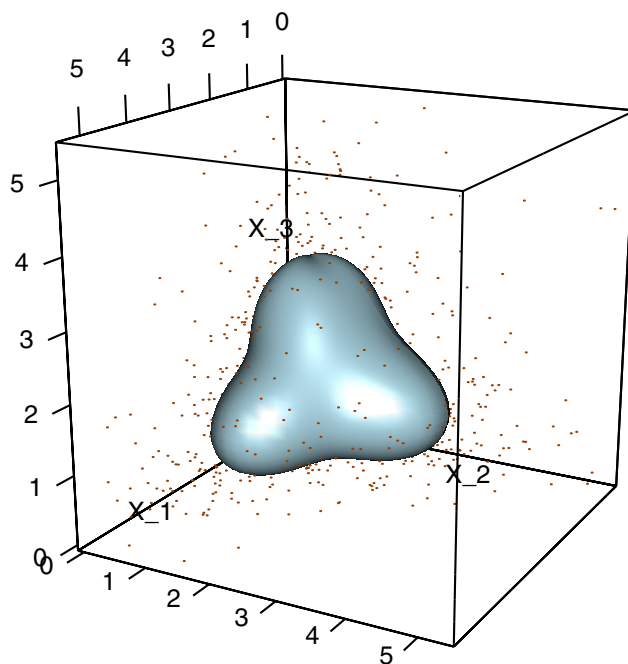
Pro příklady byl simulován náhodný výběr z trojrozměrného rozdělení, jehož složky jsou navzájem nezávislé a mají exponenciální rozdělení s parametrem 1. Rozsah výběru byl volen 2000.

Obr. 3.13 – 3.15 jsou zobrazeny odhady pro nijak neupravený výběr. Na posledním obrázku je dobře patrné výrazné zvlnění v oblastech blízko nejhlubšího bodu. To je způsobené mechanismy popsány v části 3.3.

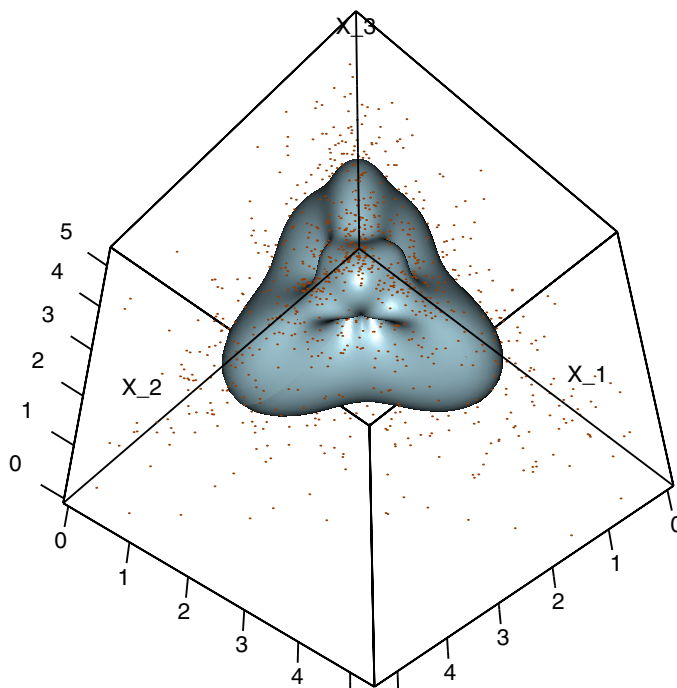
Pro odhad na obr. 3.16 byl výběr nejdříve otočen tak, že směr osy  $x_3$  přechází na směr daný spojnicí nejhlubšího bodu a počátku kartézských souřadnic. Oproti předchozímu odhadu jsou patrné výrazné hrboly právě v tomto směru (odpovídá  $\varphi_2 = \pi$ ). Na druhou stranu se odhad výrazně zlepšil v oblastech blízkých nejhlubšímu bodu. Obr. 3.17 ukazuje odhad po otočení a pak následném „zúžení“. Toto „zúžení“ je určeno maticí se sloupci  $(1, 0.05, 0.05)^T$ ,  $(0.05, 1, 0.05)^T$  a  $(0.05, 0.05, 1)^T$ .



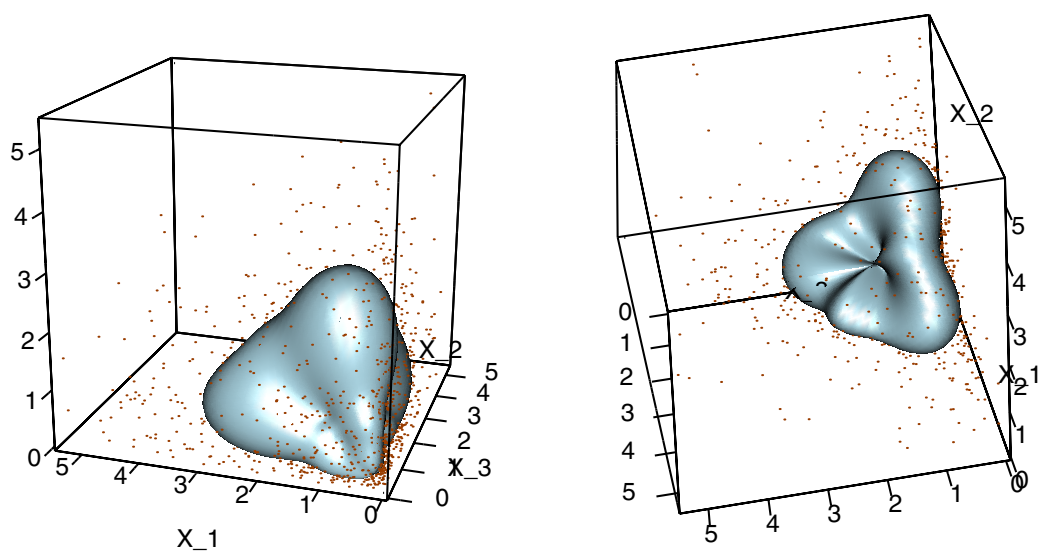
Obrázek 3.13: Výběrový směrový 0.5-quantil pro trojrozměrný výběr z exponenciálního rozdělení s parametrem 1. Rozsah výběru 2000. Řády  $p = q = 4$ .



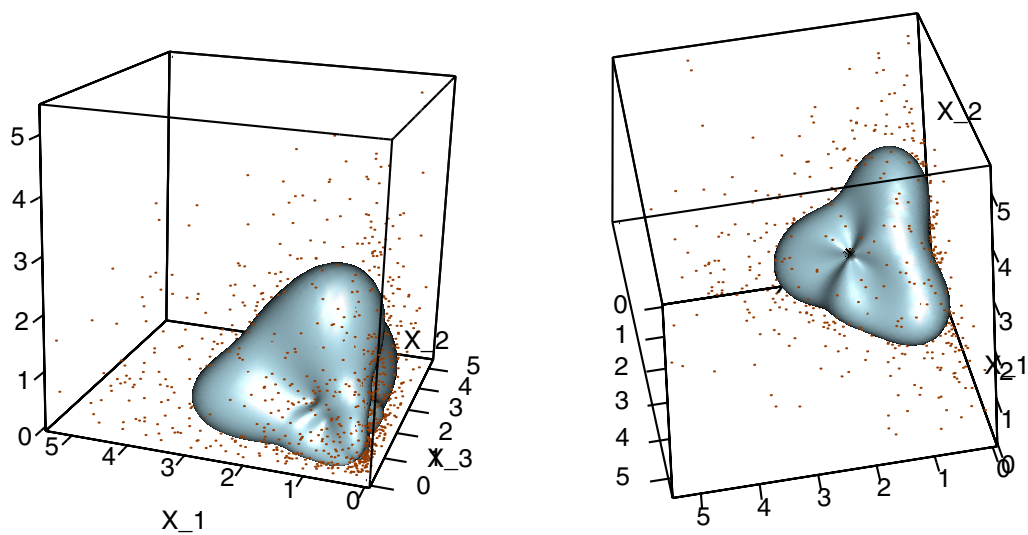
Obrázek 3.14: Výběrový 0.5-quantil pro trojrozměrný výběr z exponenciálního rozdělení s parametrem 1. Rozsah výběru 2000. Řády  $p = q = 4$ . Pohled „zepředu“.



Obrázek 3.15: Výběrový 0.5-quantil pro trojrozměrný výběr z exponenciálního rozdělení s parametrem 1. Rozsah výběru 2000. Řády  $p = q = 4$ . Pohled „zezadu“.



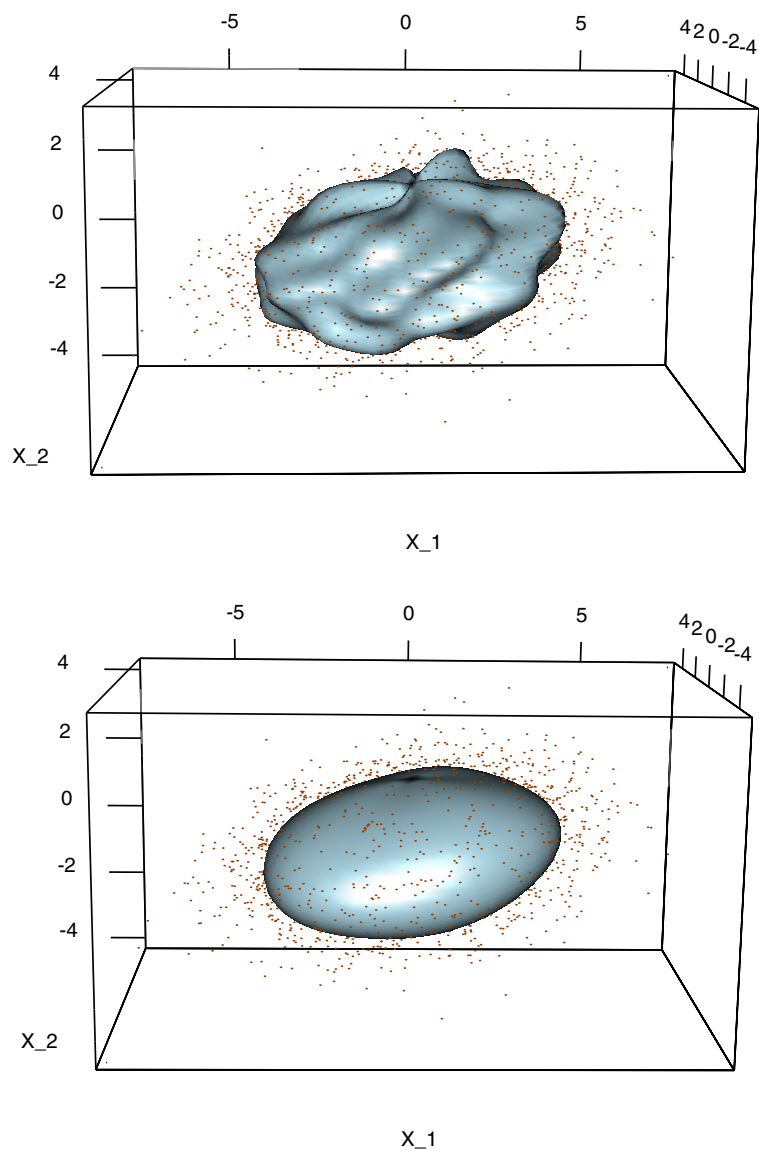
Obrázek 3.16: Výběrový kvantil odhadnutý z otočeného výběru pro řády  $p = 5$  a  $q = 4$ . Rozsah výběru 2000.



Obrázek 3.17: Výběrový kvantil odhadnutý z otočeného a zároveň „zúženého“ výběru pro řády  $p = 5$  a  $q = 4$ . Rozsah výběru 2000.



## Normální rozdělení – trojrozměrný případ



Obrázek 3.18: 75% kvantily pro výběr z normálního rozdělení s nulovou střední hodnotou a varianční maticí s diagonálními prvky 4, 2, 1. Mimodiagonální prvky jsou rovné 0.5. Rozsah výběru 5000. Podobně jako v dvojrozměrném případě jsme na data nejprve aplikovali lineární transformaci určenou maticí z Choleského rozkladu výběrové varianční matice (výběrová varianční matice transformovaného výběru je jednotková). Horní obrázek zobrazuje odhad pro řády  $p = q = 8$ . Dolní  $p = q = 3$ . Pro vyšší hodnoty  $p, q$  se tvar kvantilu nijak výrazně nemění, jen je více „hrbolatý“. V takovém případě je lepší volit menší hodnoty řádů.

### Příklad - porodní váha a délka

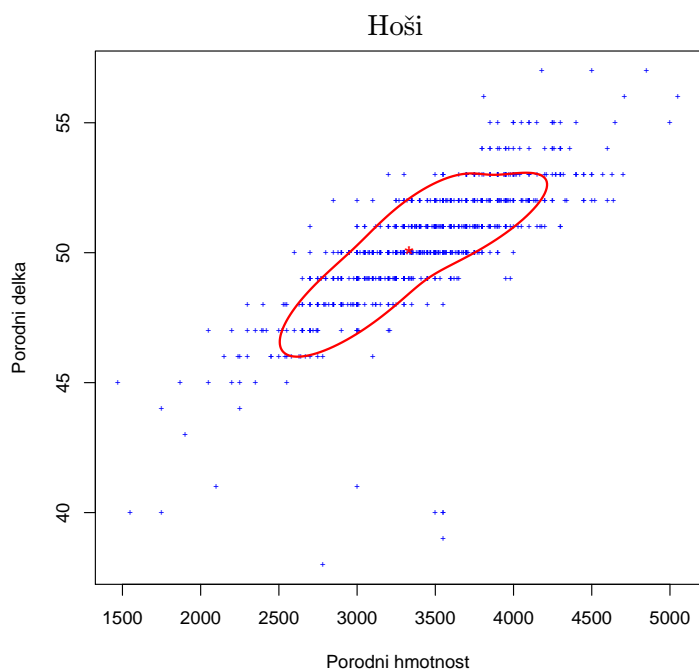
Nakonec uvedme příklad kvantilu zhotoveného pro reálná data. Zajímáme se o 75% kvantil porodní délky a hmotnosti u dětí. K dispozici máme 847 pozorování pro hochy a 786 pro dívky. Ačkoliv jsou tyto veličiny spojité, tak pro použití naší metody nemusí být zrovna nejvhodnější, jelikož porodní délka se zaokrouhlovala na celé centimetry. A vzhledem k rozmezí, ve kterém se tato veličina pohybuje, má téměř diskrétní charakter. Ale i přesto docházíme k poměrně rozumným výsledkům.

Podobně jako v uvedených příkladech na normální rozdělení přejdeme nejprve k výběru s jednotkovou varianční maticí. Po získání odhadu, pro který byl jak pro chlapce tak pro dívky volen řád  $p = 4$ , se vrátíme zpět k původnímu výběru.

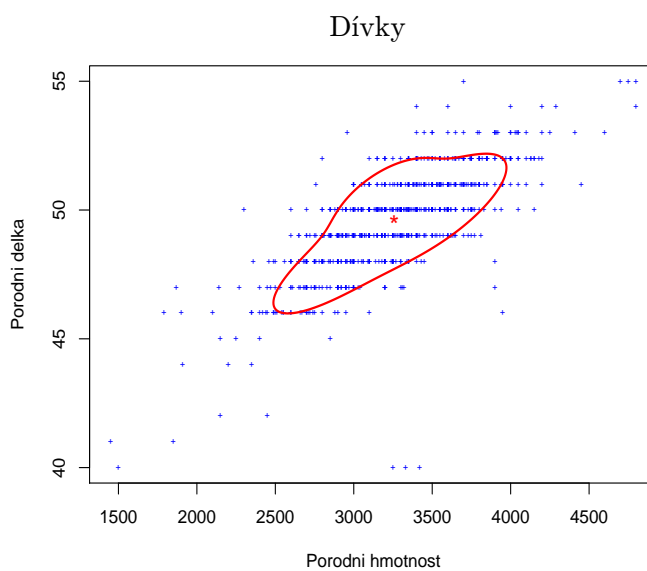
Klasickým parametrickým odhadem při předpokladech normálního rozdělení bychom dostali množinu tvaru elipsy. A odhad by navíc byl ovlivněn několika vzdálenějšími pozorováními. Navíc data mají tu vlastnost, že s rostoucí porodní délkou roste i výběrový rozptyl veličiny udávající porodní hmotnost. Pro taková data by eliptický tvar konfidenční množiny nebyl příliš vhodný. Na obr. 3.20 je vidět, že náš odhad si s touto vlastností poměrně dobře poradí.

Zaměříme se na chvíli na maximální a minimální hodnotu souřadnice pro porodní délku, které ještě leží v periodickém kvantilu (tedy budeme se zajímat o rozmezí kvantilu odpovídající této veličině). Označme je  $x_{max}$  a  $x_{min}$ . Tyto hodnoty budou vždy velmi blízko celým číslům. Je to způsobené diskrétním charakterem veličiny udávající porodní délku a tím, že kvantil vždy prochází určitým počtem bodů (viz věta 19).  $x_{max}$  a  $x_{min}$  leží v blízkosti dvou těchto bodů a jejich hodnota je vlastně těmito body určena. A tak rozmezí kvantilu ve veličině porodní délka bude vždy vymezeno přibližně celými čísly, i když teoreticky by mohlo být určeno libovolnými čísly. Nicméně pro tímto způsobem měřená a zaznamenávaná data tato vlastnost není ničemu na škodu.

Nakonec se podívejme ještě na jednu nepříjemnost, která je opět způsobena diskrétním charakterem porodní délky. Levá dolní část kvantilu u hochů je jakoby mírně vychýlena napravo. Kvantil odhadujeme ve směrech od nejhlubšího bodu (červená hvězdička), tedy jakoby na přímkách procházejících nejhlubším bodem. Přičemž je důležitý počet bodů ležících v těsné blízkosti takových přímek. Vzdálenost bodů od sebe není až tak podstatná. Diskrétní charakter se projeví v „řádkovitém“ uspořádání dat a platí, že co řádek, to jeden bod ležící v blízkosti přímků. A např. jihozápadozápadním směrem od nejhlubšího bodu protneme méně „řádků“ než třeba směrem určeným spojnicí nejhlubšího bodu a bodu s porodní hmotností 3000 g a délkou 41 cm. Označíme-li  $l$  počet řádků, který nějaký směr protíná, pak 75% kvantil bude ležet v blízkosti průsečíku s  $[0.75l]$ -tým řádkem ve směru od nejhlubšího bodu. Proto na odhad mají poměrně velký vliv vzdálenější pozorování, která nabývají hodnot porodní délky kolem 40 cm. Jejich vliv na odhad ale není způsobený jejich odlehlostí, nýbrž tím, že více z nich vždy leží v nějakém směru od nejhlubšího bodu. Tyto body navíc datům dávají až téměř „nehvězdicovitý charakter“. Pro data „neřádkovitého“ (a tedy spojitého) charakteru by k tomuto vychýlení zřejmě nedošlo.



Obrázek 3.19: 75% kvantil pro porodní délku a hmotnost u hochů. Hodnoty naměřeny u 847 hochů. Řád volen 4.



Obrázek 3.20: 75% kvantil pro porodní délku a hmotnost u dívek. Hodnoty naměřeny u 786 dívek. Řád volen 4.

# Literatura

- Bantli, F. & Hallin, M. (1999),  $L_1$ -estimation in linear models with heterogenous white noise, *Statistics & Probability Letters* **45**, 305–315.
- Donoho, D. L. & Gasko, M. (1992), Breakdown properties of location estimates based on halfspace depth and projected outlyingness, *The Annals of Statistics* **20**(4), 1803–1827.
- Koenker, R. (2005), *Quantile Regression*, Cambridge University Press.
- Kufner, A. & Kadlec, J. (1969), *Fourierovy řady*, Academia.
- Liu, R. Y. (1990), On a notation of data depth based on random simplices, *The Annals of Statistics* **18**(1), 405–414.
- Liu, R. Y., Jesse, P. M. & Singh, K. (1999), Multivariate analysis by data depth: Descriptive statistics, graphics and interference, *The Annals of Statistics* **27**(3), 783–858.
- Milasevic, P. & Ducharme, G. R. (1987), Uniqueness of spatial median, *The Annals of Statistics* **15**(3), 1332–1333.
- Oberhofer, W. (1982), The consistency of nonlinear regression minimizing the  $L_1$ -norm, *The Annals of Statistics* **10**(1), 316–319.
- Struyf, A. & Rousseeuw, P. J. (2000), High-dimensional computation of the deepest location, *Computation Statistics and Data Analysis* **34**, 415–426.