

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

BRNO UNIVERSITY OF TECHNOLOGY



FAKULTA STROJNÍHO INŽENÝRSTVÍ
ÚSTAV MATEMATIKY

FACULTY OF MECHANICAL ENGINEERING
INSTITUTE OF MATHEMATICS

METODA BOOTSTRAP A JEJÍ APLIKACE

BOOTSTRAP METHOD AND ITS APPLICATION

DIPLOMOVÁ PRÁCE

MASTER'S THESIS

AUTOR PRÁCE

AUTHOR

LUCIE PAVLÍČKOVÁ

VEDOUCÍ PRÁCE

SUPERVISOR

doc. RNDr. ZDENĚK KARPÍŠEK, CSc.

BRNO 2009

Vysoké učení technické v Brně, Fakulta strojního inženýrství

Ústav matematiky

Akademický rok: 2009/2010

ZADÁNÍ DIPLOMOVÉ PRÁCE

student(ka): Lucie Pavlíčková

který/která studuje v **magisterském studijním programu**

obor: **Matematické inženýrství (3901T021)**

Ředitel ústavu Vám v souladu se zákonem č.111/1998 o vysokých školách a se Studijním a zkušebním řádem VUT v Brně určuje následující téma diplomové práce:

Metoda bootstrap a její aplikace

v anglickém jazyce:

Bootstrap Method and its Application

Stručná charakteristika problematiky úkolu:

Popis metody bootstrap, jejích vlastností a užití na reálných datových souborech.

Cíle diplomové práce:

Princip a vlastnosti metody bootstrap. Odhady momentových a kvantilových charakteristik rozdělení pravděpodobnosti. PC realizace metody ve statistických softwarech. Aplikace bootstrapu na reálných datových souborech.

Seznam odborné literatury:

1. Chernick, M. R. Bootstrap Methods: A Practitioner's Guide. New York: Wiley, 1999.
2. Davison, A. C. and Hinkley, D. V. Bootstrap Methods and Their Application. Cambridge, England: Cambridge University Press, 1997.
3. Efron, B. and Tibshirani, R. J. An Introduction to the Bootstrap. Boca Raton, FL: CRC Press, 1994.
4. Mooney, C. Z. and Duval, R. D. Bootstrapping: A Nonparametric Approach to Statistical Inference. Sage, 1993.
5. Odborné články dle pokynů vedoucího diplomové práce.

Vedoucí diplomové práce: doc. RNDr. Zdeněk Karpíšek, CSc.

Termín odevzdání diplomové práce je stanoven časovým plánem akademického roku 2009/2010.

V Brně, dne

L.S.

prof. RNDr. Josef Šlapal, CSc.
Ředitel ústavu

doc. RNDr. Miroslav Doupovec, CSc.
Děkan fakulty

Abstrakt

Diplomová práce popisuje metodu bootstrap a její použití pro určení přesnosti odhadu, tvorbu konfidenčních intervalů a testování statistických hypotéz. Dále je předložena metoda odhadu diskrétního rozdělení pravděpodobnosti kategoriální veličiny využívající gradient kvazinormy tohoto rozdělení. Metoda bootstrap je v konkrétních příkladech aplikována k získání konfidenčního intervalu pravděpodobnostní funkce kategoriální veličiny.

Diplomová práce je součástí řešení projektu MŠMT České republiky čís. 1M06047 „Centrum pro jakost a spolehlivost výroby“, projektu Grantové agentury České republiky reg. čís. 103/08/1658 „Pokročilá optimalizace návrhu složených betonových konstrukcí“ a výzkumného záměru MŠMT České republiky čís. MSM0021630519 „Progresivní spolehlivé a trvanlivé nosné stavební konstrukce“.

Summary

The diploma thesis describes the bootstrap method and its applications in the estimate accuracy statement, in the confidence intervals generation and in the testing of statistical hypotheses. Further the method of the discrete probability estimation of the categorical quantity is presented, making use the gradient of the quasi-norm hereof distribution. On concrete examples the bootstrap method is applied in the confidence intervals forming of the categorical quantity probability function.

The diploma thesis was supported by the project of MŠMT of the Czech Republic no. 1M06047 “Centre for Quality and Reliability of Production”, by the grant of Grant Agency of the Czech Republic (Czech Science Foundation) reg. no. 103/08/1658 “Advanced optimum design of composed concrete structures” and by the research plan of MŠMT of the Czech Republic no. MSM0021630519 “Progressive reliable and durable structures”.

Klíčová slova

bootstrap, odhad parametru, přesnost odhadu, konfidenční interval, BCA, test statistické hypotézy, f-divergence, kvazinorma, diskrétní rozdělení pravděpodobnosti, gradientní odhad

Keywords

bootstrap, parameter estimate, accuracy of the estimate, confidence interval, BCA, statistical hypothesis testing, f-divergence, quasi-norm, discrete probability distribution, gradient estimate

PAVLÍČKOVÁ, L. *Metoda bootstrap a její aplikace*. Brno: Vysoké učení technické v Brně, Fakulta strojního inženýrství, 2009. 64 s. Vedoucí diplomové práce doc. RNDr. Zdeněk Karpíšek, CSc.

Prohlašuji, že jsem diplomovou práci vypracovala samostatně a že jsem uvedla všechny využitě prameny a literaturu.

Děkuji vedoucímu práce za jeho trpělivost, ochotu, vstřícnost a pomocnou ruku. Děkuji svému otci za cenné rady při psaní této práce a za její důkladné přečtení. A konečně děkuji svému manželovi za jeho obrovskou podporu, pomoc a lásku, která mi dodávala sílu v probdělých nocích.

Obsah

1	Úvod	13
2	Základní pojmy	14
2.1	Pravděpodobnost, náhodná veličina, náhodný vektor a jejich charakteristiky	14
2.2	Náhodný výběr a jeho charakteristiky	17
2.3	Rozdělení pravděpodobnosti pro aplikace	18
2.4	Odhady parametrů a testování hypotéz	20
2.5	Lineární regresní analýza	21
3	Základní metoda bootstrap — přesnost odhadu	22
3.1	Střední kvadratická chyba a tolerance chyby odhadu	22
3.2	Bootstrapový odhad střední kvadratické chyby, rozptylu, směrodatné odchylky a vychýlení	23
3.3	Parametrický bootstrap	24
4	Intervalové odhady parametrů rozdělení pravděpodobnosti	25
4.1	Pivotové odhady	25
4.1.1	Intervalový odhad střední hodnoty	25
4.1.2	Intervalový odhad rozptylu a směrodatné odchylky	26
4.1.3	Obecný pivotový konfidenční interval	28
4.2	Kvantilové odhady	29
4.2.1	Jednoduchý kvantilový konfidenční interval	29
4.2.2	Reziduový kvantilový konfidenční interval	29
4.2.3	BCA kvantilový konfidenční interval	30
4.3	Testování hypotéz	32
5	Bootstrapový výběr z časové řady	33
6	Vícerozměrný bootstrapový výběr	34
6.1	Vícerozměrný bootstrapový výběr z k -tic	34
6.2	Vícerozměrný bootstrapový výběr z odchylek	35
7	Testy hypotéz o středních hodnotách	36
7.1	Náhodný výběr z dvojrozměrného náhodného vektoru	36
7.1.1	Shodná rozdělení pravděpodobnosti odchylek	36
7.1.2	Různá rozdělení pravděpodobnosti odchylek	38

7.1.3	Testování hypotéz	39
7.2	Náhodný výběr z k -rozměrného náhodného vektoru	39
7.2.1	Shodná rozdělení pravděpodobnosti odchylek	40
7.2.2	Různá rozdělení pravděpodobnosti odchylek	41
8	Testy hypotéz o parametrech lineárního regresního modelu	43
8.1	Konfidenční interval pro β_j	43
8.2	Test hypotézy $C\beta = 0$	44
8.3	Další metody	46
9	Odhad diskrétního rozdělení pravděpodobnosti kategoriální veličiny pomocí gradientu kvazinormy	48
10	Intervalové odhady diskrétního rozdělení pravděpodobnosti kategoriální veličiny	51
10.1	Falešná kostka	51
10.2	Volební model	56
11	Závěr	63
	Použité zdroje	64

Kapitola 1

Úvod

Statistická teorie se pokouší odpovědět na tři základní otázky: jak získat data, jak tato data analyzovat a shrnout a jak ověřit jejich přesnost. Bootstrap je metoda, která přináší snadno pochopitelnou a proveditelnou odpověď na třetí otázku. Jeho princip poprvé popsal Bradley Efron, profesor Stanfordské univerzity, v roce 1979. Jednalo se tehdy o jednu z prvních metod, která ve statistice nahrazovala tradiční algebraické výpočty počítačovými simulacemi na pozorovaných datech. Bootstrap přinesl možnost odhadnout přesnost libovolného odhadu libovolného parametru. Přitom spočívá v prosté myšlence mnohonásobného opakování jednoduchého algoritmu. Navíc není závislý na centrální limitní větě, a proto ho lze s úspěchem použít i pro výběry s malým rozsahem. S rozvojem a zrychlováním počítačů se otevřely dveře pro další aplikace bootstrapu, zejména pro konstruování konfidenčních intervalů, testování statistických hypotéz a v oblasti regresní analýzy. Dnes je bootstrap stále častěji používanou metodou s širokým využitím, bylo o něm napsáno množství knih a má své pevné místo mezi matematickým softwarem.

V této diplomové práci popíšeme, co je to bootstrap a jak ho použít pro výpočet střední kvadratické chyby odhadu, rozptylu nebo vychýlení. Předložíme několik přístupů, kterými lze získat bootstrapový konfidenční interval, a ukážeme, jak s jeho pomocí testovat statistické hypotézy. Načrtneme také, jak aplikovat bootstrap na časové řady, vícerozměrné náhodné výběry a jak ho lze použít v regresní analýze.

V deváté kapitole seznámíme čtenáře s novou metodou odhadu pravděpodobnostní funkce kategoriální veličiny pomocí gradientu kvazinormy jejího rozdělení. V desáté kapitole se nakonec budeme věnovat aplikaci bootstrapu na tuto metodu a sestrojíme intervalové odhady pravděpodobnostních funkcí konkrétních kategoriálních veličin. Otestujeme při tom nový a stále se vyvíjející software Shine bootstrap.

Kapitola 2

Základní pojmy

Na úvod připomeneme některé základní pojmy z teorie pravděpodobnosti a statistiky, jako jsou náhodná veličina a náhodný vektor a jejich charakteristiky, náhodný výběr, rozdělení pravděpodobnosti, odhady parametrů, testování hypotéz nebo regresní analýza. Viz také [5].

2.1. Pravděpodobnost, náhodná veličina, náhodný vektor a jejich charakteristiky

Pokusem rozumíme realizaci určitého systému podmínek, které jsou opakovatelné a neměnné. Výsledkem pokusu je náhodný jev. Jednotlivým možným výsledkům pokusu odpovídají *elementární náhodné jevy*, které vyjadřujeme pomocí jednoprvkových množin $\{\omega\}$. Všechny výsledky pokusu tvoří množinu Ω , kterou nazýváme *základní prostor*. Platí $\omega \in \Omega$. *Náhodným jevem* A pak rozumíme podmnožinu základního prostoru, tedy $A \subseteq \Omega$. *Opačný náhodný jev* k jevu A je jev \bar{A} , který nastane právě tehdy, když nenastane jev A ; tj. \bar{A} je doplněk A v Ω . *Jevové pole* Σ na Ω je množina náhodných jevů s vlastnostmi:

1. $\Omega \in \Sigma$,
2. pro každé $A \in \Sigma$ je $\bar{A} \in \Sigma$,
3. pro každou posloupnost náhodných jevů $A_i \in \Sigma$, $i = 1, 2, \dots$, je

$$\bigcap_{i=1}^{\infty} A_i \in \Sigma.$$

Pravděpodobnost $P(A)$ náhodného jevu $A \in \Sigma$ je reálná funkce definovaná na jevovém poli Σ s vlastnostmi:

1. $P(A) \geq 0$ pro všechny náhodné jevy $A \in \Sigma$,
2. $P(\Omega) = 1$,
3. pro každou posloupnost disjunktních náhodných jevů $A_i \in \Sigma$, $i = 1, 2, \dots$, je

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i).$$

Uspořádaná trojice (Ω, Σ, P) se nazývá *pravděpodobnostní prostor*.

Náhodná veličina (náhodná proměnná) X je funkce, která nabývá náhodně reálných číselných hodnot x . Formálně jde tedy o zobrazení $X: \Omega \rightarrow \mathbb{R}$. Přitom požadujeme, aby pro každé $c \in \mathbb{R}$ množina $\{X < c\} = \{\omega \in \Omega : X(\omega) \in (-\infty, c)\}$ byla jevem, tj. $\{X < c\} \in \Sigma$ (tzv. borelovsky měřitelné zobrazení). Množina všech hodnot náhodné veličiny X se nazývá základní soubor nebo také populace. Distribuční funkce náhodné veličiny X je reálná funkce

$$F(x) = P(X < x) = P(X \in (-\infty; x)),$$

definovaná pro všechna $x \in (-\infty; \infty)$. Distribuční funkcí je náhodná veličina plně popsána a říkáme, že je dáno její rozdělení pravděpodobnosti.

Náhodná veličina X je *diskrétní* (má *diskrétní rozdělení pravděpodobnosti*), jestliže nabývá s nenulovou pravděpodobností nejvýše spočetně mnoha hodnot x_1, x_2, \dots . Její *pravděpodobnostní funkce* je posloupnost

$$p(x) = P(X = x) > 0 \text{ pro } x = x_1, x_2, \dots$$

Náhodná veličina X je *spojitá* (má *spojité rozdělení pravděpodobnosti*), jestliže má tzv. *absolutně spojitou* distribuční funkci $F(x)$; tzn. že existuje nezáporná funkce $f(x)$ taková, že pro každé $x \in (-\infty; \infty)$ je

$$F(x) = \int_{-\infty}^x f(t) dt.$$

Funkci $f(x)$ nazýváme její *hustotou pravděpodobnosti*.

Číselné charakteristiky náhodné veličiny X jsou reálná čísla, která vyjadřují její důležité vlastnosti. Nejprve uvedeme *momentové charakteristiky*.

Polohu rozdělení pravděpodobnosti náhodné veličiny X charakterizuje její *střední hodnota*

$$E(X) = \sum_x xp(x)$$

pro diskrétní náhodnou veličinu X (pokud řada konverguje absolutně),

$$E(X) = \int_{-\infty}^{\infty} xf(x) dx$$

pro spojitou náhodnou veličinu X (pokud integrál konverguje absolutně).

Míru kolísání hodnot náhodné veličiny X kolem její střední hodnoty $E(X)$ vyjadřuje její *rozptyl (variance)*

$$D(X) = E([X - E(X)]^2).$$

Odmocninu z rozptylu pak nazýváme *směrodatnou odchylkou* náhodné veličiny X a značíme $\sigma(X) = \sqrt{D(X)}$. Odtud pak plyne také alternativní označení pro rozptyl $\sigma^2(X)$.

Míru asymetrie rozdělení náhodné veličiny X vzhledem k její střední hodnotě vyjadřuje *koefficient šikmosti (asymetrie)*

$$A(X) = \frac{E([X - E(X)]^3)}{[\sigma(X)]^3}.$$

Koeficient špičatosti (excesu)

$$\gamma_2(X) = \frac{E[X - E(X)]^4}{[D(x)]^2} - 3$$

porovnává rozdělení pravděpodobnosti náhodné veličiny X s normálním rozdělením.

Nejdůležitější *kvantilovou číselnou charakteristikou* náhodné veličiny X je p -kvantil náhodné veličiny $x_p = \inf\{x; F(x) \geq p\}$, kde $0 < p < 1$. Pro spojitou náhodnou veličinu X s rostoucí distribuční funkcí je $F(x_p) = p$. Kvantil $x_{0,5}$ je *medián* náhodné veličiny X a charakterizuje polohu jejího rozdělení pravděpodobnosti.

Modus \hat{x} náhodné veličiny X je taková hodnota, v níž nabývá její pravděpodobnostní funkce nebo hustota pravděpodobnosti maximum, popř. supremum.

Uspořádanou n -tici náhodných veličin X_1, \dots, X_n nazýváme *n -rozměrným náhodným vektorem* (X_1, \dots, X_n) . Dále *simultánní (sdružená) distribuční funkce* náhodného vektoru (X_1, \dots, X_n) je reálná funkce

$$F(x_1, \dots, x_n) = P(X_1 < x_1, \dots, X_n < x_n) = P((X_1, \dots, X_n) \in (-\infty; x_1) \times \dots \times (-\infty; x_n))$$

definovaná pro všechny n -tice $(x_1, \dots, x_n) \in (-\infty; \infty)^n$. Simultánní distribuční funkcí je náhodný vektor (X_1, \dots, X_n) jednoznačně popsán a říkáme, že je dáno jeho *simultánní rozdělení pravděpodobnosti*.

Náhodný vektor je *diskrétní*, jestliže všechny jeho složky jsou diskrétní. Náhodný vektor je *spojitý*, jestliže všechny jeho složky jsou nezávislé a spojité. Jestliže u náhodného vektoru vynecháme některé jeho složky, dostaneme *marginální rozdělení pravděpodobnosti*. *Podmínné rozdělení pravděpodobnosti* dostaneme, pokud budeme uvažovat podmínku, že některé ze složek náhodného vektoru nabývají libovolné pevné hodnoty.

Mezi nejdůležitější číselné charakteristiky náhodného vektoru patří *střed (centrum)*, což je uspořádaná n -tice $(E(X_1), \dots, E(X_n))$, jejímiž složkami jsou střední hodnoty složek náhodného vektoru.

Vzájemný vztah dvou složek X_i, X_j náhodného vektoru (X_1, \dots, X_n) vyjadřuje jejich *kovariance*

$$\text{cov}(X_i, X_j) = E([x_i - E(X_i)][x_j - E(X_j)]) = E(X_i X_j) - E(X_i)E(X_j).$$

Z kovariancí se sestavuje *kovarianční matice*

$$\text{cov}(X_1, \dots, X_n) = \begin{pmatrix} D(X_1) & \text{cov}(X_1, X_2) & \cdots & \text{cov}(X_1, X_n) \\ \text{cov}(X_2, X_1) & D(X_2) & \cdots & \text{cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_n, X_1) & \text{cov}(X_n, X_2) & \cdots & D(X_n) \end{pmatrix}.$$

Mírou lineární závislosti dvou složek X_i, X_j náhodného vektoru (X_1, \dots, X_n) je jejich *koeficient korelace*

$$\rho(X_i, X_j) = \text{cov} \left(\frac{X_i - E(X_i)}{\sigma(X_i)}, \frac{X_j - E(X_j)}{\sigma(X_j)} \right) = \frac{\text{cov}(X_i, X_j)}{\sigma(X_i)\sigma(X_j)}.$$

Z koeficientů korelace se sestavuje *korelační matice* náhodného vektoru (X_1, \dots, X_n)

$$\rho(X_1, \dots, X_n) = \begin{pmatrix} 1 & \rho(X_1, X_2) & \cdots & \rho(X_1, X_n) \\ \rho(X_2, X_1) & 1 & \cdots & \rho(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \rho(X_n, X_1) & \rho(X_n, X_2) & \cdots & 1 \end{pmatrix}.$$

2.2. Náhodný výběr a jeho charakteristiky

Opakujeme-li n -krát nezávisle pokus, jehož výsledkem je hodnota náhodné veličiny X s distribuční funkcí $F(x, \theta)$, kde θ je reálný parametr (případně vektor parametrů nebo jejich funkce) daného rozdělení pravděpodobnosti, pozorujeme vlastně náhodný vektor $\mathbf{X} = (X_1, \dots, X_n)$ a předpokládáme, že jeho složky jsou nezávislé náhodné veličiny X_i se stejnou distribuční funkcí, jako má pozorovaná náhodná veličina X . Náhodný vektor $\mathbf{X} = (X_1, \dots, X_n)$ pak nazýváme *náhodným výběrem* z náhodné veličiny X nebo z jejího rozdělení pravděpodobnosti a číslo n *rozsahem náhodného výběru*. Množinu všech uvažovaných hodnot parametru θ nazýváme *parametrickým prostorem*.

Číselný vektor $\mathbf{x} = (x_1, \dots, x_n)$, který získáme při realizaci náhodného výběru, kde x_i je pozorovaná hodnota složky X_i , $i = 1, \dots, n$, nazýváme *pozorovanou hodnotou náhodného výběru* $\mathbf{X} = (X_1, \dots, X_n)$.

Funkci náhodného výběru $T(X_1, \dots, X_n)$ nazýváme *výběrovou charakteristikou* nebo *statistikou*, její hodnotu na pozorované hodnotě náhodného výběru $t = T(x_1, \dots, x_n)$ nazýváme *empirickou charakteristikou* nebo *pozorovanou hodnotou statistiky* T .

Nejvýznamnějšími výběrovými charakteristikami jsou:

1. *výběrový průměr*

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i,$$

2. *výběrový rozptyl*

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2,$$

3. *výběrová směrodatná odchylka*

$$S = \sqrt{S^2},$$

4. *výběrový koeficient korelace*

$$R = \frac{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{S(X)S(Y)}$$

pro náhodný výběr z náhodného vektoru (X, Y) , kde $S(X)$, $S(Y)$ jsou výběrové směrodatné odchylky náhodných veličin X , Y .

Základní vlastnosti výběrového průměru a výběrového rozptylu jsou:

1. Jestliže pozorovaná náhodná veličina X má střední hodnotu $E(X)$, pak

$$E(\bar{X}) = E(X).$$

2. Jestliže pozorovaná náhodná veličina X má rozptyl $D(X)$, pak

$$D(\bar{X}) = \frac{D(X)}{n}, \quad \sigma(\bar{X}) = \frac{\sigma(X)}{\sqrt{n}}, \quad E(S^2) = D(X).$$

Hodnoty empirických charakteristik jsou náhodné, při opakovaných realizacích náhodného výběru se náhodně mění. Z předcházejícího však plyne, že např. pro $n \rightarrow \infty$ rozptyl výběrového průměru $D(\bar{X}) \rightarrow 0$, takže pro dostatečně velké n je „takřka jistě“ aritmetický průměr \bar{x} blízký neznámé střední hodnotě rozdělení $E(X)$.

2.3. Rozdělení pravděpodobnosti pro aplikace

Uvedeme si některá známá rozdělení pravděpodobnosti.

Normální (Gaussovo) rozdělení $N(\mu, \sigma^2)$, kde $\mu, \sigma^2 \in \mathbb{R}$, $\sigma > 0$, má hustotu pravděpodobnosti

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], \quad x \in (-\infty, \infty),$$

střední hodnotu $E(X) = \mu$, rozptyl $D(X) = \sigma^2$ a koeficient šikmosti $A(X) = 0$.

Transformací náhodné veličiny X s normálním rozdělením $N(\mu, \sigma^2)$ na náhodnou veličinu

$$Z = \frac{X - \mu}{\sigma}$$

dostaneme *normované normální rozdělení* $N(0, 1)$. Pro jeho kvantily platí $z_{1-\alpha} = -z_\alpha$, $\alpha \in (0, 1)$.

Pearsonovo rozdělení $\chi^2(k)$ s k stupni volnosti, kde $k \in \mathbb{N}$, má hustotu pravděpodobnosti

$$f(x) = \begin{cases} \frac{e^{-x/2} x^{k/2-1}}{2^{k/2} \Gamma(\frac{k}{2})} & \text{pro } x \in (0, \infty), \\ 0 & \text{pro } x \in (-\infty, 0), \end{cases}$$

kde

$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt, \quad z > 0,$$

je tzv. *gama funkce*, střední hodnotu $E(X) = k$, rozptyl $D(X) = 2k$ a koeficient šikmosti $A(X) = 4/\sqrt{2k}$.

Studentovo rozdělení $S(k)$ s k stupni volnosti, kde $k \in \mathbb{N}$, má hustotu pravděpodobnosti

$$f(x) = \frac{\Gamma(\frac{k+1}{2})}{\sqrt{\pi k} \Gamma(\frac{k}{2})} \left(1 + \frac{x^2}{k}\right)^{-\frac{k+1}{2}}, \quad x \in (-\infty, \infty),$$

střední hodnotu $E(X) = 0$ pro $k > 1$, rozptyl $D(X) = k/(k-2)$ pro $k > 2$ a koeficient šikmosti $A(X) = 0$ pro $k > 3$. Pro jeho kvantily platí $t_{1-\alpha} = -t_\alpha$, $\alpha \in (0, 1)$.

Fisherovo-Snedecorovo rozdělení $F(k_1, k_2)$ s k_1, k_2 stupni volnosti, kde $k_1, k_2 \in \mathbb{N}$, má hustotu pravděpodobnosti

$$f(x) = \begin{cases} \frac{\left(\frac{k_1}{k_2}\right)^{\frac{k_1}{2}} x^{\frac{k_1}{2}-1} \left(1 + \frac{k_1}{k_2}x\right)^{-\frac{k_1+k_2}{2}}}{B\left(\frac{k_1}{2}, \frac{k_2}{2}\right)} & \text{pro } x \in \langle 0, \infty \rangle, \\ 0 & \text{pro } x \in (-\infty, 0), \end{cases}$$

kde

$$B(z_1, z_2) = \int_0^1 t^{z_1-1} (1-t)^{z_2-1} dt = \frac{\Gamma(z_1)\Gamma(z_2)}{\Gamma(z_1+z_2)}, \quad z_1 > 0, z_2 > 0,$$

je tzv. *beta funkce*, střední hodnotu $E(X) = \frac{k_2}{k_2-2}$ pro $k_2 > 2$ a rozptyl $D(X) = \frac{2k_2^2(k_1+k_2-2)}{k_1(k_2-2)^2(k_2-4)}$ pro $k_2 > 4$. Pro jeho kvantily platí $F_{1-\alpha}(k_1, k_2) = 1/F_\alpha(k_2, k_1)$, $\alpha \in (0, 1)$.

Multinomické rozdělení $M(n, p_1, \dots, p_m)$, kde $n \in \mathbb{N}$, $p_1, \dots, p_m > 0$, $\sum_{i=1}^m p_i < 1$, má sdruženou pravděpodobnostní funkci

$$p(x_1, x_2, \dots, x_m) = \begin{cases} \frac{n! p_1^{x_1} p_2^{x_2} \dots p_m^{x_m} \left(1 - \sum_{i=1}^m p_i\right)^{n - \sum_{i=1}^m x_i}}{x_1! x_2! \dots x_m! \left(n - \sum_{i=1}^m x_i\right)!} & \text{pro } x_1, \dots, x_m = 1, \dots, n, \sum_{i=1}^m x_i \leq n \\ 0 & \text{jinak,} \end{cases}$$

střed $(E(X_1), \dots, E(X_m)) = (np_1, \dots, np_m)$ a kovarianční matici

$$\mathbf{cov}(X_1, \dots, X_m) = \begin{pmatrix} np_1(1-p_1) & -np_1p_2 & \dots & -np_1p_m \\ -np_2p_1 & np_2(1-p_2) & \dots & -np_2p_m \\ \vdots & & \ddots & \vdots \\ -np_m p_1 & -np_m p_2 & \dots & np_m(1-p_m) \end{pmatrix}.$$

Pro jediné $m = 1$ se $M(n, p)$ nazývá *binomické rozdělení* $Bi(n, p)$.

Weibullovo rozdělení $W(b, c, \delta)$, kde $b, \delta > 0$, $c \in \mathbb{R}$, má hustotu pravděpodobnosti

$$f(x) = \frac{b}{\delta} \left(\frac{x-c}{\delta}\right)^{b-1} \exp\left[-\left(\frac{x-c}{\delta}\right)^b\right], \quad x \in \langle c, \infty \rangle,$$

střední hodnotu $E(X) = c + \delta\Gamma(1/b + 1)$, rozptyl $D(X) = \delta^2 [\Gamma(2/b + 1) - \Gamma^2(1/b + 1)]$ a koeficient šikmosti

$$A(X) = \frac{\Gamma(3/b + 1) - 3\Gamma(1/b + 1)\Gamma(2/b + 1) + 2\Gamma^3(1/b + 1)}{[\Gamma(2/b + 1) - \Gamma^2(1/b + 1)]^{3/2}}.$$

Přitom b je *parametr tvaru*, c je *prahový parametr* a δ je *parametr měřítka*. Pro $b \approx 3,6$ je Weibullovo rozdělení blízké normálnímu rozdělení.

Pro $b = 1$ se $W(1, c, \delta)$ nazývá *exponenciální rozdělení* $E(\lambda, c)$.

2.4. Odhady parametrů a testování hypotéz

Odhadem T parametru θ je statistika $T(X_1, \dots, X_n)$, která na celém parametrickém prostoru nabývá hodnot blízkých parametru θ . Odhad T parametru θ je *nestranný* (*nevychýlený*), jestliže $E(T) = \theta$. V opačném případě je odhad *stranný* (*vychýlený*).

Bodovým odhadem parametru θ je pozorovaná hodnota $t = T(x_1, \dots, x_n)$ odhadu T . *Interval spolehlivosti* nebo také *konfidenční interval* pro parametr θ se spolehlivostí $1 - \alpha$, kde $\alpha \in \langle 0; 1 \rangle$, je dvojice takových statistik T_1, T_2 , že $P(T_1 \leq \theta \leq T_2) = 1 - \alpha$ pro libovolnou hodnotu parametru θ . *Intervalový odhad parametru θ* se spolehlivostí $1 - \alpha$ je pak interval $\langle t_1; t_2 \rangle$, kde t_1, t_2 jsou pozorované hodnoty statistik T_1, T_2 . Intervalové odhady dělíme na *oboustranné* a *jednostranné* podle toho, zda je ohraničujeme oboustranně nebo jednostranně.

Statistická hypotéza H je tvrzení o vlastnostech rozdělení pravděpodobnosti pozorované náhodné veličiny X s distribuční funkcí $F(x, \theta)$. Platnost hypotézy ověřujeme postupem, který se nazývá *test statistické hypotézy*. Proti testované (nulové) hypotéze H stavíme tzv. *alternativní hypotézu \bar{H}* , kterou volíme dle požadavků úlohy. Hypotéza je *jednoduchá*, jestliže uvažujeme jedinou hypotetickou hodnotu, v opačném případě je hypotéza *složená*. Složená hypotéza může být *jednostranná* (např. $H: \theta < \theta_0$) nebo *oboustranná* (např. $H: \theta \neq \theta_0$). *Parametrická hypotéza* je tvrzení o parametrech pozorované náhodné veličiny, *neparametrická hypotéza* je tvrzení o kvalitativních vlastnostech pozorované náhodné veličiny.

Testové kritérium je vhodná statistika $T(X_1, \dots, X_n)$ zkonstruovaná pro test dané hypotézy H proti dané alternativní hypotéze \bar{H} . Obor hodnot testového kritéria T se za předpokladu platnosti hypotézy H rozdělí na dvě disjunktní podmnožiny, a to *kritický obor W_α* a jeho doplněk \bar{W}_α . Kritický obor W_α se vzhledem k alternativní hypotéze stanoví tak, aby pravděpodobnost, že testové kritérium T nabyde hodnotu z kritického oboru, byla nejvýše α . Číslo $\alpha > 0$ nazýváme *hladina významnosti testu*. O hypotéze rozhodujeme na základě *pozorované hodnoty testového kritéria $t = T(x_1, \dots, x_n)$* . Jestliže $t \in W_\alpha$, zamítáme hypotézu H a současně nezamítáme alternativní hypotézu \bar{H} na hladině významnosti α . Jestliže naopak $t \in \bar{W}_\alpha$, nezamítáme hypotézu H a současně zamítáme alternativní hypotézu \bar{H} na hladině významnosti α .

Chyba prvního druhu nastane, jestliže zamítneme platnou hypotézu. Její pravděpodobnost je hladina významnosti testu α . *Chyba druhého druhu nastane*, jestliže nezamítneme neplatnou hypotézu. Pravděpodobnost této chyby značíme β a číslo $1 - \beta$ nazýváme *silou testu*.

K testování statistických hypotéz lze rovněž použít přímo intervalové odhady. Při testování na hladině významnosti α pak místo testového kritéria zvolíme vhodný intervalový odhad se spolehlivostí $1 - \alpha$.

Při testování statistických hypotéz lze také místo kritického oboru použít tzv. *p-hodnotu*. Pro pozorovanou hodnotu t testového kritéria T je *p-hodnotou* číslo $1 - P(-t \leq T \leq t)$. Jestliže $p < \alpha$, zamítáme hypotézu H a současně nezamítáme alternativní hypotézu \bar{H} na hladině významnosti α . Jestliže naopak $p \geq \alpha$, nezamítáme hypotézu H a současně zamítáme alternativní hypotézu \bar{H} na hladině významnosti α .

2.5. Lineární regresní analýza

Regresní analýza zkoumá závislost závisle proměnné náhodné veličiny Y na nezávisle proměnném náhodném vektoru $\mathbf{X} = (X_1, \dots, X_k)$. *Základní lineární regresní model* uvažujeme ve tvaru

$$Y_i = \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \varepsilon_i, \quad i = 1, \dots, n,$$

vektorově zapisujeme

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

kde

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} X_{11} & \cdots & X_{k1} \\ \vdots & \ddots & \vdots \\ X_{1n} & \cdots & X_{kn} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_k \end{pmatrix},$$

přičemž $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ a $\text{cov}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$, kde $\sigma^2 > 0$, takže složky vektoru $\boldsymbol{\varepsilon}$ jsou nekorelované a mají stejný rozptyl σ^2 . Platí

$$E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}, \quad \text{cov}(\mathbf{Y}) = \sigma^2 \mathbf{I}.$$

Nejlepší nestranný odhad vektoru *regresních koeficientů* $\boldsymbol{\beta}$ je vektor \mathbf{b} , který získáme *metodou nejmenších čtverců*, tj. minimalizací *reziduálního součtu čtverců*

$$(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}).$$

Jestliže matice \mathbf{X} má *plnou hodnotu* $k < n$, existuje jediný vektor \mathbf{b} , který je řešením *soustavy normálních rovnic*

$$\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{Y},$$

kde matice $\mathbf{X}^T \mathbf{X}$ je regulární, a to

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

Platí

$$E(\mathbf{b}) = \boldsymbol{\beta}, \quad \text{cov}(\mathbf{b}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}.$$

Vektor bodových odhadů \hat{Y}_i hodnot (středních hodnot) Y_i , $i = 1, \dots, n$, je

$$\hat{\mathbf{Y}} = \mathbf{X}\mathbf{b} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

Bodový odhad rozptylu σ^2 je

$$s^2 = \frac{\mathbf{Y}^T \mathbf{Y} - \mathbf{b}^T \mathbf{X}^T \mathbf{Y}}{n - k} = \frac{\mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}}{n - k} = \frac{\mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \hat{\mathbf{Y}}}{n - k}. \quad (2.1)$$

Kapitola 3

Základní metoda bootstrap — přesnost odhadu

3.1. Střední kvadratická chyba a tolerance chyby odhadu

Nechť X je náhodná veličina a nechť θ je parametr jejího rozdělení pravděpodobnosti. Realizujeme náhodný výběr (X_1, \dots, X_n) z této náhodné veličiny o rozsahu n . Na základě pozorovaných hodnot vypočítáme odhad parametru θ a označíme jej $\hat{\theta}$. Pak *střední kvadratickou chybou odhadu* $\hat{\theta}$ rozumíme hodnotu

$$\text{MSE} = E(\hat{\theta} - \theta)^2.$$

Důležitost střední kvadratické chyby vyplývá z následující *Čebyševovy-Markovovy nerovnosti*:

$$P(|\hat{\theta} - \theta| \leq k\sqrt{\text{MSE}}) \geq 1 - \frac{1}{k^2} \text{ pro libovolné } k.$$

Pokud např. zvolíme $k = 2$, dostáváme

$$P(\theta - 2\sqrt{\text{MSE}} \leq \hat{\theta} \leq \theta + 2\sqrt{\text{MSE}}) \geq 0,75.$$

Číslo $2\sqrt{\text{MSE}}$ nazýváme *tolerancí chyby odhadu* a používáme ho jako hrubou míru přesnosti odhadu v případech, kdy žádnou vhodnější míru nemáme k dispozici.

Vychýlením odhadu $\hat{\theta}$ parametru θ rozumíme $\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$. Snadno se ukáže, že platí $\text{MSE} = D(\hat{\theta}) + (\text{Bias}(\hat{\theta}))^2$.

Pokud například odhadujeme střední hodnotu rozdělení $E(X)$ výběrovým průměrem \bar{X} , můžeme využít faktu, že rozptyl výběrového průměru je

$$D(\bar{X}) = \frac{D(X)}{n}.$$

Odhadneme-li dále rozptyl rozdělení $D(X)$ výběrovým rozptylem S^2 , pak s využitím centrální limitní věty dostáváme asymptotický konfidenční interval se spolehlivostí 0,95:

$$E(X) \in \left(\bar{X} - 2\frac{S}{\sqrt{n}}; \bar{X} + 2\frac{S}{\sqrt{n}} \right).$$

Pak $2Sn^{-1/2}$ je v tomto případě tolerancí chyby odhadu \bar{X} .

V mnoha případech ovšem neznáme analytickou metodu výpočtu střední kvadratické chyby a tolerance chyby odhadu. Představme si ale, že bychom mohli mnohokrát opakovat realizaci náhodného výběru o rozsahu n z náhodné veličiny X . Označme $\hat{\theta}_i$ odhad parametru θ vypočítaný z pozorovaných hodnot náhodného výběru při i -tém opakování a MSE_i střední kvadratickou chybu tohoto odhadu. Pak při dostatečně velkém počtu opakování můžeme střední kvadratickou chybu odhadu $\hat{\theta}$ odhadnout

$$\widehat{MSE} = \frac{1}{B} \sum_{i=1}^B MSE_i = \frac{1}{B} \sum_{i=1}^B (\hat{\theta}_i - \theta)^2,$$

kde číslo B označuje počet provedených opakování. Takovéto opakování obvykle není ve skutečnosti možné, a proto musíme přistoupit k další aproximaci výpočtu odhadu MSE.

3.2. Bootstrapový odhad střední kvadratické chyby, rozptylu, směrodatné odchylky a vychýlení

Pokud neznáme rozdělení pravděpodobnosti pozorované náhodné veličiny X nebo není k dispozici intervalový odhad jejího parametru θ , nahradíme soubor pozorovaných hodnot (x_1, \dots, x_n) náhodného výběru (X_1, \dots, X_n) novým souborem získaným z (x_1, \dots, x_n) náhodným výběrem s opakováním (s vrácením). Takto získaný náhodný výběr nazýváme *bootstrapovým výběrem*.

Při odhadu střední kvadratické chyby, rozptylu, směrodatné odchylky a vychýlení odhadu $\hat{\theta}$ parametru θ rozdělení pravděpodobnosti náhodné veličiny X postupujeme následovně:

1. Z pozorovaných hodnot (x_1, \dots, x_n) náhodného výběru (X_1, \dots, X_n) vypočítáme odhad $\hat{\theta}$ parametru θ .
2. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu n z pozorovaných hodnot (x_1, \dots, x_n) . Obvykle volíme $B \geq 1000$.
3. Pro každý bootstrapový výběr vypočítáme odhad parametru θ a označíme $\hat{\theta}_{b,i}$, kde $i = 1, 2, \dots, B$.
4. *Bootstrapovým odhadem parametru θ* rozumíme obvykle aritmetický průměr

$$\hat{\theta}_b = \frac{1}{B} \sum_{i=1}^B \hat{\theta}_{b,i}.$$

5. *Bootstrapovým odhadem střední kvadratické chyby MSE odhadu $\hat{\theta}$* je

$$\widehat{MSE}_b = \frac{1}{B} \sum_{i=1}^B (\hat{\theta}_{b,i} - \hat{\theta})^2.$$

6. *Bootstrapovým odhadem rozptylu $D(\hat{\theta})$* je

$$\widehat{D}(\hat{\theta})_b = \frac{1}{B-1} \sum_{i=1}^B \left(\hat{\theta}_{b,i} - \frac{1}{B} \sum_{i=1}^B \hat{\theta}_{b,i} \right)^2.$$

7. *Bootstrapovým odhadem směrodatné odchylky* $\sigma(\hat{\theta})$ je

$$\hat{\sigma}(\hat{\theta})_b = \sqrt{\frac{1}{B-1} \sum_{i=1}^B \left(\hat{\theta}_{b,i} - \frac{1}{B} \sum_{i=1}^B \hat{\theta}_{b,i} \right)^2}. \quad (3.1)$$

8. *Bootstrapovým odhadem vychýlení Bias*($\hat{\theta}$) odhadu $\hat{\theta}$ je

$$\hat{B}(\hat{\theta})_b = \frac{1}{B} \sum_{i=1}^B \hat{\theta}_{b,i} - \hat{\theta}.$$

3.3. Parametrický bootstrap

Na metodu bootstrap lze také pohlížet *parametricky*. V tom případě předpokládáme, že známe rozdělení pravděpodobnosti pozorované náhodné veličiny X . Z pozorovaných hodnot (x_1, \dots, x_n) náhodného výběru (X_1, \dots, X_n) odhadneme všechny potřebné parametry její distribuční funkce a bootstrapové náhodné výběry nyní realizujeme jako náhodné výběry z takto odhadnuté distribuční funkce. Všechny následné výpočty se pak už neliší od těch popsaných výše.

Kapitola 4

Intervalové odhady parametrů rozdělení pravděpodobnosti

4.1. Pivotové odhady

Nechť Z je spojitá náhodná veličina se střední hodnotou $E(Z) = 0$, rozptylem $D(Z) = 1$ a hustotou pravděpodobnosti $f(z)$. Nechť X je spojitá náhodná veličina daná vztahem

$$X = \mu + \sigma Z, \quad \text{kde } \sigma > 0,$$

tedy má hustotu pravděpodobnosti

$$g(x) = \frac{1}{\sigma} f\left(\frac{x - \mu}{\sigma}\right).$$

Potom střední hodnota $E(X) = \mu$, rozptyl $D(X) = \sigma^2$ a směrodatná odchylka $\sigma(X) = \sigma$.

V této podkapitole ukážeme, jak získat metodou bootstrap odhad konfidenčního intervalu pro odhady střední hodnoty μ , rozptylu σ^2 a směrodatné odchylky σ . Přitom budeme používat následující odhady parametrů: μ odhadneme výběrovým průměrem \bar{X} a σ odhadneme výběrovou směrodatnou odchylkou S . Více viz [3]

V následujícím budeme ukazovat konfidenční intervaly se spolehlivostí $1 - 2\alpha$ získané pomocí α -kvantilů a $(1 - \alpha)$ -kvantilů příslušných rozdělení pravděpodobnosti. Pochopitelně není nutné volit právě tyto kvantily. Konfidenční interval se spolehlivostí ξ získáme pomocí libovolné dvojice ζ -kvantil, η -kvantil, kde $\eta - \zeta = \xi$ a $\xi, \zeta, \eta \in \langle 0; 1 \rangle$.

4.1.1. Intervalový odhad střední hodnoty

Pokud Z má normované normální rozdělení pravděpodobnosti $N(0; 1)$, pak statistika

$$t = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

má Studentovo rozdělení pravděpodobnosti s $n - 1$ stupni volnosti a platí

$$P\left(-t_{1-\alpha} < \frac{\bar{X} - \mu}{S/\sqrt{n}} < t_{1-\alpha}\right) = 1 - 2\alpha,$$

kde $t_{1-\alpha}$ je $(1 - \alpha)$ -kvantil Studentova rozdělení s $n - 1$ stupni volnosti. Odtud odvodíme konfidenční interval se spolehlivostí $1 - 2\alpha$

$$\mu \in \left(\bar{X} - t_{1-\alpha} \frac{S}{\sqrt{n}}; \bar{X} + t_{1-\alpha} \frac{S}{\sqrt{n}} \right).$$

Pokud Z nemá normální rozdělení pravděpodobnosti, rozdělení pravděpodobnosti statistiky t je stále nezávislé na μ i σ , ale už se nejedná o Studentovo rozdělení. Nicméně pokud bychom mohli nějakým způsobem zjistit hodnoty kvantilů tohoto neznámého rozdělení, stále by platilo

$$P \left(t_\alpha < \frac{\bar{X} - \mu}{S/\sqrt{n}} < t_{1-\alpha} \right) = 1 - 2\alpha$$

a konfidenční interval se spolehlivostí $1 - 2\alpha$ by měl tvar

$$\mu \in \left(\bar{X} - t_{1-\alpha} \frac{S}{\sqrt{n}}; \bar{X} - t_\alpha \frac{S}{\sqrt{n}} \right).$$

Hodnoty kvantilů rozdělení pravděpodobnosti statistiky t odhadneme pomocí metody bootstrap.

Postup pro získání konfidenčního intervalu pro $\mu = E(X)$ bude tedy následující:

1. Z pozorovaných hodnot (x_1, \dots, x_n) náhodného výběru (X_1, \dots, X_n) vypočítáme pozorované hodnoty výběrového průměru \bar{X} a výběrové směrodatné odchyly S .
2. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu n z pozorovaných hodnot (x_1, \dots, x_n) . Obvykle volíme $B \geq 1000$.
3. Pro každý bootstrapový výběr vypočítáme pozorovanou hodnotu výběrového průměru $\bar{X}_{b,i}$ a výběrové směrodatné odchyly $S_{b,i}$ a hodnotu statistiky t

$$t_{b,i} = \frac{\bar{X}_{b,i} - \bar{X}}{S_{b,i}/\sqrt{n}},$$

kde $i = 1, 2, \dots, B$.

4. α -kvantil a $(1-\alpha)$ -kvantil rozdělení pravděpodobnosti statistiky t_b odhadneme hodnotami $t_{b,\alpha}$ a $t_{b,1-\alpha}$ splňujícími co nejpřesněji

$$|\{t_{b,i}; t_{b,i} \leq t_{b,\alpha}\}|/B \doteq \alpha, \quad |\{t_{b,i}; t_{b,i} \leq t_{b,1-\alpha}\}|/B \doteq 1 - \alpha,$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

5. *Bootstrapovým konfidenčním intervalem se spolehlivostí $1 - 2\alpha$ pro střední hodnotu $E(X)$ je*

$$\mu \in \left(\bar{X} - t_{b,1-\alpha} \frac{S}{\sqrt{n}}; \bar{X} - t_{b,\alpha} \frac{S}{\sqrt{n}} \right).$$

4.1.2. Intervalový odhad rozptylu a směrodatné odchyly

Pokud Z má normované normální rozdělení pravděpodobnosti $N(0; 1)$, pak statistika

$$\chi^2 = \frac{(n-1)S^2}{\sigma^2}$$

má Pearsonovo rozdělení pravděpodobnosti s $n - 1$ stupni volnosti a platí

$$P \left(\chi_{\alpha}^2 < \frac{(n-1)S^2}{\sigma^2} < \chi_{1-\alpha}^2 \right) = 1 - 2\alpha,$$

kde χ_{α}^2 a $\chi_{1-\alpha}^2$ jsou α -kvantil a $(1 - \alpha)$ -kvantil Pearsonova rozdělení s $n - 1$ stupni volnosti. Odtud odvodíme konfidenční interval se spolehlivostí $1 - 2\alpha$

$$\sigma^2 \in \left(\frac{(n-1)S^2}{\chi_{1-\alpha}^2}; \frac{(n-1)S^2}{\chi_{\alpha}^2} \right).$$

Pokud Z nemá normální rozdělení pravděpodobnosti, rozdělení pravděpodobnosti statistiky χ^2 je stále nezávislé na μ i σ , ale už se nejedná o Pearsonovo rozdělení. Nicméně pokud bychom mohli nějakým způsobem zjistit hodnoty kvantilů tohoto neznámého rozdělení, výše uvedené rovnosti by stále platily. Hodnoty kvantilů rozdělení pravděpodobnosti statistiky χ^2 odhadneme pomocí metody bootstrap.

Postup pro získání konfidenčního intervalu pro $\sigma^2 = D(X)$ a $\sigma = \sigma(X)$ bude tedy následující:

1. Z pozorovaných hodnot (x_1, \dots, x_n) náhodného výběru (X_1, \dots, X_n) vypočítáme pozorovanou hodnotu výběrového rozptylu S^2 .
2. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu n z pozorovaných hodnot (x_1, \dots, x_n) . Obvykle volíme $B \geq 1000$.
3. Pro každý bootstrapový výběr vypočítáme pozorovanou hodnotu výběrového rozptylu $S_{b,i}^2$ a hodnotu statistiky χ^2

$$\chi_{b,i}^2 = \frac{(n-1)S_{b,i}^2}{S^2},$$

kde $i = 1, 2, \dots, B$.

4. α -kvantil a $(1 - \alpha)$ -kvantil rozdělení pravděpodobnosti statistiky χ_b^2 odhadneme hodnotami $\chi_{b,\alpha}^2$ a $\chi_{b,1-\alpha}^2$ splňujícími co nejpřesněji

$$|\{\chi_{b,i}^2; \chi_{b,i}^2 \leq \chi_{b,\alpha}^2\}|/B \doteq \alpha, \quad |\{\chi_{b,i}^2; \chi_{b,i}^2 \leq \chi_{b,1-\alpha}^2\}|/B \doteq 1 - \alpha,$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

5. *Bootstrapovým konfidenčním intervalem se spolehlivostí $1 - 2\alpha$ pro rozptyl $D(X)$ je*

$$\sigma^2 \in \left(\frac{(n-1)S^2}{\chi_{b,1-\alpha}^2}; \frac{(n-1)S^2}{\chi_{b,\alpha}^2} \right).$$

6. *Bootstrapovým konfidenčním intervalem se spolehlivostí $1 - 2\alpha$ pro směrodatnou odchylku $\sigma(X)$ je*

$$\sigma \in \left(\sqrt{\frac{(n-1)S^2}{\chi_{b,1-\alpha}^2}}; \sqrt{\frac{(n-1)S^2}{\chi_{b,\alpha}^2}} \right).$$

4.1.3. Obecný pivotový konfidenční interval

Výše uvedené konfidenční intervaly jsou konkrétními příklady obecnějších formulí. Zcela obecné konfidenční intervaly pro libovolné parametry rozdělení pravděpodobnosti odhadované libovolnými statistikami uvádí [2].

Nechť θ je libovolný parametr (popř. parametrická funkce) rozdělení pravděpodobnosti náhodné veličiny X a nechť $\hat{\theta}$ je nějakým jeho odhadem. Nechť $\hat{\sigma}(\hat{\theta})_b$ je bootstrapovým odhadem směrodatné odchyly odhadu $\hat{\theta}$ podle (3.1). Pokud u náhodné veličiny $\hat{\theta}_b$, tj. odhadu parametru θ na základě bootstrapových náhodných výběrů, můžeme předpokládat normální rozdělení pravděpodobnosti, pak statistika

$$t = \frac{\hat{\theta} - \theta}{\hat{\sigma}(\hat{\theta})_b}$$

má Studentovo rozdělení pravděpodobnosti s $n - 1$ stupni volnosti a platí

$$P\left(-t_{1-\alpha} < \frac{\hat{\theta} - \theta}{\hat{\sigma}(\hat{\theta})_b} < t_{1-\alpha}\right) = 1 - 2\alpha,$$

kde $t_{1-\alpha}$ je $(1 - \alpha)$ -kvantil Studentova rozdělení pravděpodobnosti s $n - 1$ stupni volnosti. Odtud snadno odvodíme bootstrapový konfidenční interval se spolehlivostí $1 - 2\alpha$ pro odhad parametru θ

$$\theta \in (\hat{\theta} - t_{1-\alpha}\hat{\sigma}(\hat{\theta})_b; \hat{\theta} + t_{1-\alpha}\hat{\sigma}(\hat{\theta})_b).$$

V případě, kdy u náhodné veličiny $\hat{\theta}_b$ nelze předpokládat normální rozdělení, musíme hodnoty kvantilů rozdělení pravděpodobnosti statistiky t opět odhadnout metodou bootstrap.

Postup pro získání konfidenčního intervalu pro parametr θ bude tedy následující:

1. Z pozorovaných hodnot (x_1, \dots, x_n) náhodného výběru (X_1, \dots, X_n) vypočítáme odhad $\hat{\theta}$ parametru θ .
2. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu n z pozorovaných hodnot (x_1, \dots, x_n) . Obvykle volíme $B \geq 1000$.
3. Pro každý bootstrapový výběr vypočítáme odhad $\hat{\theta}_{b,i}$ parametru θ a jeho směrodatnou odchytku $\sigma(\hat{\theta}_{b,i})$. Pokud pro ni neznáme žádné analytické vyjádření, odhadneme ji vnořenou metodou bootstrap podle (3.1). Obvykle volíme $B \geq 100$.
4. Pro každý bootstrapový výběr dále vypočítáme hodnotu statistiky

$$t_{b,i} = \frac{\hat{\theta}_{b,i} - \hat{\theta}}{\sigma(\hat{\theta}_{b,i})},$$

kde $i = 1, 2, \dots, B$.

5. Podle (3.1) spočítáme odhad směrodatné odchyly $\hat{\sigma}(\hat{\theta})_b$ odhadu $\hat{\theta}$.
6. α -kvantil a $(1 - \alpha)$ -kvantil rozdělení pravděpodobnosti statistiky t_b odhadneme hodnotami $t_{b,\alpha}$ a $t_{b,1-\alpha}$ splňujícími co nejpřesněji

$$|\{t_{b,i}; t_{b,i} \leq t_{b,\alpha}\}|/B \doteq \alpha, \quad |\{t_{b,i}; t_{b,i} \leq t_{b,1-\alpha}\}|/B \doteq 1 - \alpha,$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

7. Bootstrapovým konfidenčním intervalem se spolehlivostí $1 - 2\alpha$ pro parametr θ je

$$\theta \in (\hat{\theta} - t_{b,1-\alpha}\hat{\sigma}(\hat{\theta})_b; \hat{\theta} - t_{b,\alpha}\hat{\sigma}(\hat{\theta})_b).$$

4.2. Kvantilové odhady

V této podkapitole se seznámíme s intervalovými odhady, které vycházejí přímo z rozdělení pravděpodobnosti bodových odhadů. Jsou proto zcela obecné, použitelné pro libovolný parametr, příp. parametrickou funkci, a pro libovolný jeho odhad.

4.2.1. Jednoduchý kvantilový konfidenční interval

Postup pro získání jednoduchého kvantilového konfidenčního intervalu je následující:

1. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu n z pozorovaných hodnot (x_1, \dots, x_n) náhodného výběru (X_1, \dots, X_n) . Obvykle volíme $B \geq 1000$.
2. Pro každý bootstrapový výběr vypočítáme odhad $\hat{\theta}_{b,i}$ parametru θ .
3. α -kvantil a $(1-\alpha)$ -kvantil rozdělení pravděpodobnosti statistiky $\hat{\theta}_b$ odhadneme hodnotami $\hat{\theta}_{b,\alpha}$ a $\hat{\theta}_{b,1-\alpha}$ splňujícími co nejpřesněji

$$|\{\hat{\theta}_{b,i}; \hat{\theta}_{b,i} \leq \hat{\theta}_{b,\alpha}\}|/B \doteq \alpha, \quad |\{\hat{\theta}_{b,i}; \hat{\theta}_{b,i} \leq \hat{\theta}_{b,1-\alpha}\}|/B \doteq 1 - \alpha,$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

4. *Bootstrapovým jednoduchým kvantilovým konfidenčním intervalem se spolehlivostí $1-2\alpha$ pro parametr θ je*

$$\theta \in (\hat{\theta}_{b,\alpha}; \hat{\theta}_{b,1-\alpha}).$$

4.2.2. Reziduový kvantilový konfidenční interval

Reziduem rozumíme $\varepsilon = \hat{\theta} - \theta$. Označíme α -kvantil a $(1-\alpha)$ -kvantil rozdělení pravděpodobnosti náhodné veličiny ε jako ε_α a $\varepsilon_{1-\alpha}$. Pak platí

$$P(\varepsilon_\alpha < \hat{\theta} - \theta \leq \varepsilon_{1-\alpha}) = 1 - 2\alpha.$$

Odtud odvodíme konfidenční interval se spolehlivostí $1 - 2\alpha$

$$\theta \in (\hat{\theta} - \varepsilon_{1-\alpha}; \hat{\theta} - \varepsilon_\alpha).$$

Kvantily rozdělení pravděpodobnosti rezidua ovšem neznáme a odhadneme je metodou bootstrap.

Postup pro získání reziduového kvantilového konfidenčního intervalu je následující:

1. Z pozorovaných hodnot (x_1, \dots, x_n) náhodného výběru (X_1, \dots, X_n) vypočítáme odhad $\hat{\theta}$ parametru θ .
2. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu n z pozorovaných hodnot (x_1, \dots, x_n) . Obvykle volíme $B \geq 1000$.
3. Pro každý bootstrapový výběr vypočítáme odhad $\hat{\theta}_{b,i}$ parametru θ a reziduum $e_{b,i} = \hat{\theta}_{b,i} - \hat{\theta}$.
4. α -kvantil a $(1-\alpha)$ -kvantil rozdělení pravděpodobnosti reziduí e_b odhadneme hodnotami $e_{b,\alpha}$ a $e_{b,1-\alpha}$ splňujícími co nejpřesněji

$$|\{e_{b,i}; e_{b,i} \leq e_{b,\alpha}\}|/B \doteq \alpha, \quad |\{e_{b,i}; e_{b,i} \leq e_{b,1-\alpha}\}|/B \doteq 1 - \alpha,$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

5. *Bootstrapovým reziduovým kvantilovým konfidenčním intervalem se spolehlivostí $1 - 2\alpha$ pro parametr θ je*

$$\theta \in \langle \widehat{\theta} - e_{b,1-\alpha}; \widehat{\theta} - e_{b,\alpha} \rangle.$$

4.2.3. BCA kvantilový konfidenční interval

Může se stát, že jednoduché a reziduové kvantilové konfidenční intervaly jsou vychýlené nebo příliš široké oproti hodnotám pozorovaným v praxi. K odstranění těchto nedostatků proto konstruujeme tzv. *BCA konfidenční intervaly* (z anglického bias corrected and accelerated), které jsou sice také omezeny dvěma kvantily rozdělení pravděpodobnosti bootstrapového odhadu $\widehat{\theta}_b$, ale na rozdíl od předchozích metod se již nemusí nutně jednat o α -kvantil a $(1 - \alpha)$ -kvantil pro spolehlivost $1 - 2\alpha$.

V této metodě vycházíme z předpokladu, že existuje nějaká transformace parametru θ , jejíž rozdělení pravděpodobnosti je normální a jejíž střední hodnota a rozptyl závisí na θ . Konfidenční interval pak zkonstruujeme pro transformovaný parametr a inverzní transformací jeho mezí získáme konfidenční interval pro θ . Elegance metody spočívá v tom, že předpokládanou transformaci vůbec nepotřebujeme znát v explicitním vyjádření, realizujeme ji metodou bootstrap.

Předpokládejme, že existuje rostoucí transformační zobrazení T takové, že $T(\widehat{\theta})$ má normální rozdělení pravděpodobnosti se střední hodnotou

$$E[T(\widehat{\theta})] = T(\theta) - z_0[1 + aT(\theta)]$$

a směrodatnou odchylkou

$$\sigma[T(\widehat{\theta})] = 1 + aT(\theta).$$

Nechť $z_{1-\alpha}$ je $(1 - \alpha)$ -kvantil normovaného normálního rozdělení pravděpodobnosti. Pak

$$P\left(-z_{1-\alpha} < \frac{T(\widehat{\theta}) - T(\theta)}{1 + aT(\theta)} + z_0 < z_{1-\alpha}\right) = 1 - 2\alpha,$$

odkud snadno odvodíme konfidenční interval se spolehlivostí $1 - 2\alpha$

$$T(\theta) \in \left(\frac{T(\widehat{\theta}) + z_0 - z_{1-\alpha}}{1 - a(z_0 - z_{1-\alpha})}; \frac{T(\widehat{\theta}) + z_0 + z_{1-\alpha}}{1 - a(z_0 + z_{1-\alpha})}\right).$$

Vzhledem k tomu, že náhodné veličiny

$$\frac{T(\widehat{\theta}_b) - T(\widehat{\theta})}{1 + aT(\widehat{\theta})} + z_0 \quad \text{a} \quad \frac{T(\widehat{\theta}) - T(\theta)}{1 + aT(\theta)} + z_0$$

mají stejné rozdělení pravděpodobnosti (podle předpokladu normované normální), platí

$$\begin{aligned} P\left(T(\widehat{\theta}_b) < \frac{T(\widehat{\theta}) + z_0 + z_{1-\alpha}}{1 - a(z_0 + z_{1-\alpha})}\right) &= P\left(\frac{T(\widehat{\theta}_b) - T(\widehat{\theta})}{1 + aT(\widehat{\theta})} + z_0 < \frac{z_0 + z_{1-\alpha}}{1 - a(z_0 + z_{1-\alpha})} + z_0\right) = \\ &= P\left(Z < \frac{z_0 + z_{1-\alpha}}{1 - a(z_0 + z_{1-\alpha})} + z_0\right), \end{aligned}$$

kde Z má normované normální rozdělení pravděpodobnosti. Odtud vyplývá, že horní mez konfidenčního intervalu pro $T(\theta)$ je

$$z_H = \frac{z_0 + z_{1-\alpha}}{1 - a(z_0 + z_{1-\alpha})} + z_0.$$

Obdobně se odvodí dolní mez

$$z_D = \frac{z_0 - z_{1-\alpha}}{1 - a(z_0 - z_{1-\alpha})} + z_0.$$

Zbývá odhadnout hodnoty z_0 a a . Nechť p_0 je podíl pozorovaných hodnot rozdělení pravděpodobnosti náhodné veličiny $\hat{\theta}_b$ s vlastností $\hat{\theta}_b \leq \hat{\theta}$ ku všem pozorovaným hodnotám. Pak z_0 má takovou hodnotu, že platí $P(Z \leq z_0) = p_0$, kde Z má normované normální rozdělení pravděpodobnosti. Z toho vyplývá, že z_0 koriguje vychýlení mediánu bootstrapového odhadu.

Akcelerace a měří rychlost změny $\hat{\sigma}(\hat{\theta})_b$ (viz rovnice (3.1)) v závislosti na změně skutečné hodnoty parametru θ . Uvádí se více různých odhadů pro a , nejčastěji však následující založený na míře špičatosti rozdělení náhodné veličiny X .

Označme $\hat{\theta}_{-i}$ odhad parametru θ spočítaný s vynecháním X_i , tj. z náhodného výběru $(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$. Označme

$$\hat{\theta}_{(\cdot)} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{-i}.$$

Pak

$$a = \frac{\sum_{i=1}^n (\hat{\theta}_{(\cdot)} - \hat{\theta}_{-i})^3}{\sigma \left[\sum_{i=1}^n (\hat{\theta}_{(\cdot)} - \hat{\theta}_{-i})^2 \right]^{3/2}}. \quad (4.1)$$

Postup pro získání BCA konfidenčního intervalu pro parametr θ je tedy následující:

1. Z pozorovaných hodnot (x_1, \dots, x_n) náhodného výběru (X_1, \dots, X_n) vypočítáme odhad $\hat{\theta}$ parametru θ .
2. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu n z pozorovaných hodnot (x_1, \dots, x_n) . Obvykle volíme $B \geq 1000$.
3. Pro každý bootstrapový výběr vypočítáme odhad $\hat{\theta}_{b,i}$ parametru θ .
4. Spočítáme korekci vychýlení mediánu

$$z_0 = \Phi^{-1} \left(\frac{|\{\hat{\theta}_{b,i}; \hat{\theta}_{b,i} < \hat{\theta}\}|}{B} \right),$$

kde Φ je distribuční funkce normovaného normálního rozdělení pravděpodobnosti a $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

5. Spočítáme akceleraci a podle vzorce (4.1).
6. Spočítáme

$$\alpha_1 = \Phi \left(\frac{z_0 - z_{1-\alpha}}{1 - a(z_0 - z_{1-\alpha})} + z_0 \right) \quad \text{a} \quad \alpha_2 = \Phi \left(\frac{z_0 + z_{1-\alpha}}{1 - a(z_0 + z_{1-\alpha})} + z_0 \right).$$

7. α_1 -kvantil a $(1 - \alpha_2)$ -kvantil rozdělení pravděpodobnosti statistiky $\widehat{\theta}_b$ odhadneme hodnotami $\widehat{\theta}_{b,\alpha_1}$ a $\widehat{\theta}_{b,1-\alpha_2}$ splňujícími co nejpřesněji

$$\left| \{ \widehat{\theta}_{b,i}; \widehat{\theta}_{b,i} \leq \widehat{\theta}_{b,\alpha_1} \} \right| / B \doteq \alpha_1, \quad \left| \{ \widehat{\theta}_{b,i}; \widehat{\theta}_{b,i} \leq \widehat{\theta}_{b,1-\alpha_2} \} \right| / B \doteq 1 - \alpha_2.$$

8. *BCA konfidenčním intervalem se spolehlivostí $1 - 2\alpha$ pro parametr θ je*

$$\theta \in (\widehat{\theta}_{b,\alpha_1}; \widehat{\theta}_{b,1-\alpha_2}).$$

Jednodušší verzí BCA intervalů jsou tzv. *BC konfidenční intervaly*, u kterých se volí $a = 0$.

4.3. Testování hypotéz

Bootstrapové konfidenční intervaly mají zásadní využití také v oblasti testování hypotéz. Chceme-li testovat hypotézu $H: \theta = \theta_0$ na hladině významnosti 2α , sestojíme jednoduše konfidenční interval pro θ se spolehlivostí $1 - 2\alpha$. Hypotézu na hladině významnosti 2α nezamítneme, pokud θ_0 bude prvkem intervalu, a zamítneme, pokud θ_0 nebude prvkem intervalu.

Kvantily ke konstrukci konfidenčního intervalu přitom zvolíme intuitivně v závislosti na tvaru alternativy. Testujeme-li např. proti alternativě $\overline{H}: \theta \neq \theta_0$, budeme při konstrukci konfidenčního intervalu libovolného typu volit α -kvantil a $(1 - \alpha)$ -kvantil. Testujeme-li proti alternativě $\overline{H}: \theta < \theta_0$, budeme při konstrukci reziduového kvantilového a pivotových konfidenčních intervalů volit 2α -kvantil a 1-kvantil, zatímco při konstrukci ostatních kvantilových konfidenčních intervalů $(1 - 2\alpha)$ -kvantil a 0-kvantil. Přesně naopak pak při testování proti alternativě $\overline{H}: \theta > \theta_0$, tj. při konstrukci reziduového kvantilového a pivotových konfidenčních intervalů $(1 - 2\alpha)$ -kvantil a 0-kvantil, zatímco při konstrukci ostatních kvantilových konfidenčních intervalů 2α -kvantil a 1-kvantil.

Kapitola 5

Bootstrapový výběr z časové řady

Nechť (Ω, Σ, P) je pravděpodobnostní prostor a $T \subset \mathbb{R}$. Náhodným (stochastickým) procesem nazveme reálnou funkci $X(\omega, t)$ definovanou na množině $\Omega \times T$ takovou, že při každém pevném $t_0 \in T$ je $X(\omega, t_0)$ náhodnou veličinou. Stručně značíme náhodný proces pouze $\{X(t)\}_{t \in T}$. Jestliže množina parametrů T je diskrétní (spočetná), je $\{X(t)\}_{t \in T}$ náhodný proces s diskrétním časem. Nazýváme ho též náhodnou posloupností nebo časovou řadou. Je-li množina parametrů T nespočetná, je $\{X(t)\}_{t \in T}$ náhodný proces se spojitým časem. Nazýváme ho též náhodnou funkcí.

Náhodný proces $\{X(t)\}_{t \in T}$ nazýváme Markovův (markovský), jestliže pro libovolná reálná čísla a, b a pro libovolná $t_1 < \dots < t_n < t$, kde $t_1, \dots, t_n, t \in T$ platí tzv. markovská vlastnost

$$P(a \leq X(t) < b | X(t_1) = x_1, \dots, X(t_n) = x_n) = P(a \leq X(t) < b | X(t_n) = x_n).$$

Chceme-li metodu bootstrap aplikovat na Markovův proces, je třeba pro bootstrapový výběr zachovat markovskou vlastnost. Toho nejlépe dosáhneme vhodným rekurentním vyjádřením, odhadem odchylek a realizací bootstrapového výběru z odchylek.

Nechť (X_1, \dots, X_n) je markovská časová řada, pro kterou platí $X_i = h(X_{i-1}) + \varepsilon_i$ pro každé $i = 2, \dots, n$, kde odchylky $\varepsilon_2, \dots, \varepsilon_n$ tvoří náhodný výběr z neznámé distribuční funkce s nulovou střední hodnotou. Pak bootstrapovou časovou řadu získáme takto:

1. Odhadneme funkci $h(X)$ odhadem $\hat{h}(X)$ a pro každé $i = 2, \dots, n$ vypočítáme odhad odchylky ε_i jako $e_i = X_i - \hat{h}(X_{i-1})$.
2. Realizujeme bootstrapový náhodný výběr o rozsahu $n-1$ z hodnot e_2, \dots, e_n a označíme ho $e_{b,2}, \dots, e_{b,n}$.
3. Bootstrapová časová řada se tvoří rekurzivně: $X_{b,1} = X_1$, $X_{b,i} = \hat{h}(X_{i-1}) + e_{b,i}$ pro $i = 2, \dots, n$.
4. Dále vygenerujeme B takových bootstrapových časových řad, které použijeme k získání potřebných odhadů podle postupů uvedených dříve.

V [2] se uvádí další metoda získání bootstrapové časové řady z časové řady (X_1, \dots, X_n) :

1. Zvolíme vhodnou délku bloku d tak, že každá dvě pozorování z původní časové řady vzdálená od sebe o víc než d kroků jsou již nezávislá. Uvažujeme pak všechny bloky délky d , tj. $(d+1)$ -tice X_i, \dots, X_{i+d} . (Např. pro markovskou časovou řadu bude $d = 1$.)
2. Bootstrapovou časovou řadu získáme náhodným výběrem s opakováním z těchto bloků o rozsahu k , kde $n \doteq (d+1)k$.

Kapitola 6

Vícerozměrný bootstrapový výběr

Je známo více metod, kterými lze získat vícerozměrný bootstrapový náhodný výběr. Při jeho realizaci obvykle nemůžeme postupovat u každé složky samostatně, protože chceme zachovat případnou závislost (kovarianci, korelaci), která mezi složkami může existovat. Pouze pokud předpokládáme nezávislost složek nebo případná závislost nemá vliv na statistiky, které nás dále zajímají, můžeme si dovolit samostatný bootstrapový výběr pro každou složku. Uvedeme si dva základní přístupy k realizaci vícerozměrného bootstrapového výběru.

6.1. Vícerozměrný bootstrapový výběr z k -tic

Nechť (X_{1i}, \dots, X_{ki}) , kde $i = 1, 2, \dots, n$, je náhodný výběr o rozsahu n z k -rozměrného náhodného vektoru (X_1, \dots, X_k) . Potom k -rozměrným bootstrapovým výběrem z k -tic rozumíme náhodný výběr s opakováním z k -tic (X_{1i}, \dots, X_{ki}) , tj. $((X_{1b_1}, \dots, X_{kb_1}), \dots, (X_{1b_n}, \dots, X_{kb_n}))$, kde b_1, \dots, b_n je náhodný výběr s opakováním z čísel $1, \dots, n$. Takto získaný k -rozměrný bootstrapový náhodný výběr se používá zejména pro $k = 2$ k výpočtům konfidenčních intervalů dvojitých charakteristik, jako jsou kovariance, koeficient korelace nebo poměr středních hodnot. Je také vhodný pro lineární a zejména nelineární regresní analýzu.

Uvedeme si např. postup pro získání konfidenčního intervalu pro koeficient korelace:

1. Realizujeme B bootstrapových náhodných výběrů z pozorovaných hodnot náhodného vektoru (X, Y) . Obvykle volíme $B \geq 1000$.
2. Pro každý bootstrapový výběr spočítáme pozorovanou hodnotu výběrového koeficientu korelace $R_{b,i}$, $i = 1, \dots, B$.
3. Pro intervalový odhad koeficientu korelace $\rho(X, Y)$ použijeme jednoduchý kvantilový konfidenční interval nebo lépe BCA konfidenční interval.

Reziduový kvantilový konfidenční interval se nedoporučuje, protože v některých případech dává meze mimo interval $(-1, 1)$. Použití obecného pivotového konfidenčního intervalu je také možné, ale v tomto případě výrazně výpočtově náročnější.

6.2. Vícerozměrný bootstrapový výběr z odchylek

Předpokládáme, že závislost náhodných veličin X_1, \dots, X_k, Y můžeme vyjádřit následujícím vztahem:

$$Y = h(X_1, \dots, X_k) + \varepsilon, \quad (6.1)$$

kde $h(X_1, \dots, X_k)$ je nějaká funkce náhodných veličin X_1, \dots, X_k a *odchylka* ε je na nich nezávislá náhodná proměnná se střední hodnotou 0 a směrodatnou odchylkou σ . Myšlenkou této metody je, že zafixujeme hodnoty náhodných veličin X_1, \dots, X_k , odhadneme odchylku ε_i pro každé $i = 1, \dots, n$ a realizujeme bootstrapový výběr z odchylek. Dosazením do vztahu (6.1) dostaneme pro každou původní hodnotu náhodných veličin X_1, \dots, X_k novou hodnotu náhodné veličiny Y .

Postup pro získání dvojrozměrného bootstrapového výběru je tedy následující:

1. Vypočítáme odhad $\hat{h}(X_1, \dots, X_k)$ funkce $h(X_1, \dots, X_k)$.
2. Pro každé $i = 1, \dots, n$ vypočítáme odhad odchylky $e_i = Y_i - \hat{h}(X_{1i}, \dots, X_{ki})$.
3. Realizujeme náhodný výběr s opakováním z hodnot e_1, \dots, e_n a označíme $e_{b,1}, \dots, e_{b,n}$.
4. Pro každé $i = 1, \dots, n$ vypočítáme $Y_{b,i} = \hat{h}(X_{1i}, \dots, X_{ki}) + e_{b,i}$.
5. $(k+1)$ -rozměrný bootstrapový výběr pak dostáváme ve tvaru $((X_{11}, \dots, X_{k1}, Y_{b,1})$ až $(X_{1n}, \dots, X_{kn}, Y_{b,n}))$.

Tento přístup používáme samozřejmě tehdy, když hodnoty náhodných veličin X_1, \dots, X_k byly v praxi skutečně zafixované, např. se měřily hodnoty náhodné veličiny Y v pravidelných časových a prostorových odstupech. Je také vhodná v případech regresní analýzy, pokud X_1 až X_k chápeme jako nezávisle proměnné a Y jako závisle proměnnou.

Odhady odchylek e_1, \dots, e_n mají ale rozptyly závislé na náhodných veličinách X_1 až X_k . Označme

$$\mathbf{X} = \begin{pmatrix} X_{11} & \cdots & X_{k1} \\ \vdots & \ddots & \vdots \\ X_{1n} & \cdots & X_{kn} \end{pmatrix} \quad \text{a} \quad \mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T.$$

Definujme pro každé $i = 1, \dots, n$ i -tý diagonální prvek matice \mathbf{H} jako *vliv i -tého pozorování* Y_i a označme h_i . Pak rozptyl e_i je $\sigma^2(1 - h_i)$, kde σ^2 je rozptyl odchylky ε . Zavedeme pro každé $i = 1, \dots, n$

$$r_i = \frac{e_i}{\sqrt{1 - h_i}}.$$

Pak *korigovanou odchylkou* rozumíme pro každé $i = 1, \dots, n$ veličinu

$$e_i^* = r_i - \frac{1}{n} \sum_{i=1}^n r_i.$$

Korigované odchylky mají konstantní rozptyl σ^2 stejně jako ε a jejich součet je nulový jako součet pozorovaných hodnot e_1, \dots, e_n . Proto se někdy doporučuje realizovat bootstrapový výběr nikoliv přímo z pozorovaných hodnot e_1, \dots, e_n , ale z korigovaných odchylek e_1^* až e_n^* . Zejména je tento postup vhodný, pokud rozsah náhodného výběru n je velmi malý nebo pokud jsou mezi pozorovanými hodnotami extrémně odlehlá pozorování.

Kapitola 7

Testy hypotéz o středních hodnotách

7.1. Náhodný výběr z dvojrozměrného náhodného vektoru

Uvažujme nyní dvě náhodné veličiny Y_1, Y_2 a náhodné výběry z těchto náhodných veličin $(Y_{11}, \dots, Y_{1n_1}), (Y_{21}, \dots, Y_{2n_2})$, kde Y_{ij} označuje j -té pozorování náhodné veličiny $Y_i, j = 1, \dots, n_i, i = 1, 2$. Označme dále μ_i střední hodnotu a $F_i(y)$ distribuční funkci náhodné veličiny Y_i . Odchylky definujeme jako $\varepsilon_{ij} = Y_{ij} - \mu_i$.

Budeme se zabývat problémem, jak vytvořit konfidenční interval pro rozdíl středních hodnot $\mu_1 - \mu_2$. Pomocí tohoto konfidenčního intervalu pak také můžeme testovat hypotézu $H: \mu_1 - \mu_2 = \Delta$. Jejím speciálním případem pro $\Delta = 0$ je hypotéza $H: \mu_1 = \mu_2$. Budeme přitom uvažovat dva případy.

V prvním případě budeme předpokládat, že rozdělení pravděpodobnosti odchylek jsou pro obě náhodné veličiny shodná. Pak lze ukázat, že rozdělení pravděpodobnosti obou náhodných veličin Y_1, Y_2 jsou stejného typu a jsou pouze navzájem posunutá o $\Delta = \mu_1 - \mu_2$, tj. $F_1(y) = F_2(y - \Delta)$. Pokud odchylky mají normální rozdělení pravděpodobnosti, nastává tento případ právě tehdy, když rozptyly náhodných veličin Y_1, Y_2 jsou shodné, tj. $D(Y_1) = D(Y_2)$.

V druhém případě předpokládáme, že rozdělení pravděpodobnosti odchylek jsou různá. Pokud se přitom jedná o dvě různá normální rozdělení, nastává tento případ právě tehdy, když rozptyly náhodných veličin Y_1, Y_2 jsou různé, tj. $D(Y_1) \neq D(Y_2)$.

7.1.1. Shodná rozdělení pravděpodobnosti odchylek

Předpokládáme shodná rozdělení pravděpodobnosti odchylek pro obě náhodné veličiny Y_1, Y_2 , tj. existuje jediná náhodná veličina ε s distribuční funkcí $F(\varepsilon)$. Označme

$$t = \frac{\bar{Y}_1 - \bar{Y}_2 - (\mu_1 - \mu_2)}{S_p \sqrt{1/n_1 + 1/n_2}},$$

kde \bar{Y}_1, \bar{Y}_2 jsou příslušné výběrové průměry a

$$S_p = \sqrt{\sum_{i=1}^2 \sum_{j=1}^{n_i} \frac{(Y_{ij} - \bar{Y}_i)^2}{n_1 + n_2 - 2}}.$$

Rozdělení pravděpodobnosti statistiky t závisí pouze na $F(\varepsilon)$, což se ukáže snadno, když dosadíme $\mu_i + \varepsilon_{ij}$ za Y_{ij} . Po úpravě dostaneme

$$t_\varepsilon = \frac{\bar{\varepsilon}_1 - \bar{\varepsilon}_2}{S_{\varepsilon p} \sqrt{1/n_1 + 1/n_2}},$$

kde $\bar{\varepsilon}_1, \bar{\varepsilon}_2$ jsou výběrové průměry odchylek pro náhodné veličiny Y_1, Y_2 a

$$S_{\varepsilon p} = \sqrt{\sum_{i=1}^2 \sum_{j=1}^{n_i} \frac{(\varepsilon_{ij} - \bar{\varepsilon}_i)^2}{n_1 + n_2 - 2}}.$$

Pokud ε má normální rozdělení pravděpodobnosti, pak statistiky t i t_ε mají Studentovo rozdělení pravděpodobnosti s $n_1 + n_2 - 2$ stupni volnosti a platí

$$P\left(-t_{1-\alpha} < \frac{\bar{Y}_1 - \bar{Y}_2 - (\mu_1 - \mu_2)}{S_p \sqrt{1/n_1 + 1/n_2}} < t_{1-\alpha}\right) = 1 - 2\alpha,$$

kde $t_{1-\alpha}$ je $(1 - \alpha)$ -kvantil Studentova rozdělení. Odtud konfidenční interval se spolehlivostí $1 - 2\alpha$ pro rozdíl $\mu_1 - \mu_2$ dostaneme ve tvaru

$$\mu_1 - \mu_2 \in \left(\bar{Y}_1 - \bar{Y}_2 - t_{1-\alpha} S_p \sqrt{1/n_1 + 1/n_2}; \bar{Y}_1 - \bar{Y}_2 + t_{1-\alpha} S_p \sqrt{1/n_1 + 1/n_2}\right).$$

Pokud ε nemá normální rozdělení pravděpodobnosti, odhadneme hodnoty kvantilů rozdělení pravděpodobnosti statistiky t_ε metodou bootstrap. Bootstrapový konfidenční interval pro rozdíl $\mu_1 - \mu_2$ získáme následujícím postupem:

1. Z pozorovaných hodnot náhodných výběrů z náhodných veličin Y_1, Y_2 vypočítáme pozorované hodnoty výběrových průměrů \bar{Y}_1, \bar{Y}_2 a odchylek $e_{ij} = Y_{ij} - \bar{Y}_i$.
2. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu $n_1 + n_2$ z pozorovaných hodnot odchylek $e_{11}, \dots, e_{1n_1}, e_{21}, \dots, e_{2n_2}$ a označíme $e_{b,11}, \dots, e_{b,1n_1}, e_{b,21}, \dots, e_{b,2n_2}$. Obvykle volíme $B \geq 1000$.
3. Pro každý bootstrapový výběr vypočítáme pozorované hodnoty výběrových průměrů $\bar{e}_{b,i}$ a hodnotu statistiky $t_{e,b}$

$$t_{e,b,k} = \frac{\bar{e}_{b,1} - \bar{e}_{b,2}}{S_{b,ep} \sqrt{1/n_1 + 1/n_2}}, \quad \text{kde} \quad S_{b,ep} = \sqrt{\sum_{i=1}^2 \sum_{j=1}^{n_i} \frac{(e_{b,ij} - \bar{e}_{b,i})^2}{n_1 + n_2 - 2}}.$$

kde $k = 1, \dots, B$.

4. α -kvantil a $(1 - \alpha)$ -kvantil rozdělení pravděpodobnosti statistiky $t_{e,b}$ odhadneme hodnotami $t_{b,\alpha}$ a $t_{b,1-\alpha}$ splňujícími co nejpřesněji

$$\left|\{t_{e,b,k}; t_{e,b,k} \leq t_{b,\alpha}\}\right|/B \doteq \alpha, \quad \left|\{t_{e,b,k}; t_{e,b,k} \leq t_{b,1-\alpha}\}\right|/B \doteq 1 - \alpha,$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

5. Bootstrapovým konfidenčním intervalem se spolehlivostí $1 - 2\alpha$ pro rozdíl středních hodnot $\mu_1 - \mu_2$ je

$$\mu_1 - \mu_2 \in \left(\bar{Y}_1 - \bar{Y}_2 - t_{b,1-\alpha} S_p \sqrt{1/n_1 + 1/n_2}; \bar{Y}_1 - \bar{Y}_2 - t_{b,\alpha} S_p \sqrt{1/n_1 + 1/n_2}\right).$$

7.1.2. Různá rozdělení pravděpodobnosti odchylek

Pokud rozdělení pravděpodobnosti odchylek pro náhodné veličiny Y_1, Y_2 jsou různá, musíme změnit statistiku i způsob realizace bootstrapových výběrů z odchylek. Označme

$$z = \frac{\bar{Y}_1 - \bar{Y}_2 - (\mu_1 - \mu_2)}{\sqrt{S_1^2/n_1 + S_2^2/n_2}},$$

kde S_1^2, S_2^2 jsou výběrové rozptyly náhodných výběrů z Y_1, Y_2 . Rozdělení pravděpodobnosti statistiky z závisí pouze na rozdělení pravděpodobnosti odchylek, což se ukáže snadno, dosadíme-li $\mu_i + \varepsilon_{ij}$ za Y_{ij} . Po úpravě dostaneme

$$z_\varepsilon = \frac{\bar{\varepsilon}_1 - \bar{\varepsilon}_2}{\sqrt{S_{\varepsilon 1}^2/n_1 + S_{\varepsilon 2}^2/n_2}},$$

kde

$$S_{\varepsilon i}^2 = \sum_{j=1}^{n_i} \frac{(\varepsilon_{ij} - \bar{\varepsilon}_i)^2}{n_i - 1}, \quad i = 1, 2.$$

Hodnoty kvantilů rozdělení pravděpodobnosti statistiky z_ε odhadneme metodou bootstrap. Bootstrapový konfidenční interval pro rozdíl $\mu_1 - \mu_2$ získáme následujícím postupem:

1. Z pozorovaných hodnot náhodných výběrů z náhodných veličin Y_1, Y_2 vypočítáme pozorované hodnoty výběrových průměrů \bar{Y}_1, \bar{Y}_2 a odchylek $e_{ij} = Y_{ij} - \bar{Y}_i$.
2. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu n_1 z pozorovaných hodnot odchylek e_{11}, \dots, e_{1n_1} a o rozsahu n_2 z pozorovaných hodnot odchylek e_{21}, \dots, e_{2n_2} a označíme $e_{b,11}, \dots, e_{b,1n_1}$ a $e_{b,21}, \dots, e_{b,2n_2}$. Obvykle volíme $B \geq 1000$.
3. Pro každý bootstrapový výběr vypočítáme pozorované hodnoty výběrových průměrů $\bar{e}_{b,i}$ a hodnotu statistiky $z_{e,b}$

$$z_{e,b,k} = \frac{\bar{e}_{b,1} - \bar{e}_{b,2}}{\sqrt{S_{b,e1}^2/n_1 + S_{b,e2}^2/n_2}}, \quad \text{kde} \quad S_{b,ei}^2 = \sum_{j=1}^{n_i} \frac{(e_{b,ij} - \bar{e}_{b,i})^2}{n_i - 1}, \quad i = 1, 2,$$

pro $k = 1, 2, \dots, B$.

4. α -kvantil a $(1 - \alpha)$ -kvantil rozdělení pravděpodobnosti statistiky $z_{e,b}$ odhadneme hodnotami $z_{b,\alpha}$ a $z_{b,1-\alpha}$ splňujícími co nejpřesněji

$$|\{z_{e,b,k}; z_{e,b,k} \leq z_{b,\alpha}\}|/B \doteq \alpha, \quad |\{z_{e,b,k}; z_{e,b,k} \leq z_{b,1-\alpha}\}|/B \doteq 1 - \alpha,$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

5. Bootstrapovým konfidenčním intervalem se spolehlivostí $1 - 2\alpha$ pro rozdíl středních hodnot $\mu_1 - \mu_2$ je

$$\mu_1 - \mu_2 \in \left(\bar{Y}_1 - \bar{Y}_2 - z_{b,1-\alpha} \sqrt{S_1^2/n_1 + S_2^2/n_2}; \bar{Y}_1 - \bar{Y}_2 - z_{b,\alpha} \sqrt{S_1^2/n_1 + S_2^2/n_2} \right).$$

Další metody

Je nepříjemné, že ani v případě normálních rozdělení pravděpodobnosti odchylek nemají statistiky z a z_ε Studentovo rozdělení pravděpodobnosti a hodnoty jejich kvantilů je třeba odhadovat metodou bootstrap. To lze obejít tzv. *Satterthwaitovou aproximací*, podle které má statistika z Studentovo rozdělení s K stupni volnosti, kde

$$K \doteq \frac{(S_1^2/n_1 + S_2^2/n_2)^2}{\frac{(S_1^2/n_1)^2}{n_1 - 1} + \frac{(S_2^2/n_2)^2}{n_2 - 1}}.$$

Namísto statistiky z můžeme také použít přímo statistiku $\bar{Y}_1 - \bar{Y}_2$ a pro ni sestrojít některý z dříve uvedených kvantilových konfidenčních intervalů. V tom případě realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu n_1 přímo z pozorovaných hodnot y_{11}, \dots, y_{1n_1} a o rozsahu n_2 z pozorovaných hodnot y_{21}, \dots, y_{2n_2} .

7.1.3. Testování hypotéz

Popsali jsme konstrukci konfidenčních intervalů pivotového typu pro statistiku $\mu_1 - \mu_2$, využili jsme přitom α -kvantily a $(1 - \alpha)$ -kvantily rozdělení pravděpodobnosti odchylek. Jiné kvantily těchto rozdělení se získají analogicky.

Testujeme-li hypotézu $H: \mu_1 - \mu_2 = \Delta$ proti alternativní hypotéze $\bar{H}: \mu_1 - \mu_2 \neq \Delta$, využijeme příslušný konfidenční interval sestrojený pomocí α -kvantilu a $(1 - \alpha)$ -kvantilu. Testujeme-li ale proti alternativní hypotéze $\bar{H}: \mu_1 - \mu_2 > \Delta$, využijeme příslušný konfidenční interval sestrojený pomocí $(1 - 2\alpha)$ -kvantilu a 0-kvantilu. A konečně testujeme-li proti alternativní hypotéze $\bar{H}: \mu_1 - \mu_2 < \Delta$, využijeme příslušný konfidenční interval sestrojený pomocí 2α -kvantilu a 1-kvantilu.

Pokud pozorovaná hodnota parametru Δ je prvkem intervalového odhadu, pak na hladině významnosti 2α nezamítáme hypotézu H a zároveň zamítáme alternativní hypotézu \bar{H} . Pokud pozorovaná hodnota parametru Δ není prvkem intervalového odhadu, pak zamítáme hypotézu H a zároveň nezamítáme alternativní hypotézu \bar{H} .

7.2. Náhodný výběr z k -rozměrného náhodného vektoru

Rozšířme nyní úvahy z minulé části na obecnější případ. Uvažujme k náhodných veličin Y_1 až Y_k a náhodné výběry z těchto náhodných veličin $(Y_{11}, \dots, Y_{1n_1}), \dots, (Y_{k1}, \dots, Y_{kn_k})$, kde Y_{ij} označuje j -té pozorování náhodné veličiny Y_i , $j = 1, \dots, n_i$, $i = 1, \dots, k$. Označme dále μ_i střední hodnotu náhodné veličiny Y_i . Odchylky definujeme jako $\varepsilon_{ij} = Y_{ij} - \mu_i$.

Budeme se zabývat problémem, jak testovat hypotézu $H: \mu_1 = \dots = \mu_k = 0$ proti alternativní hypotéze $\bar{H}: \exists i (\mu_i \neq 0)$. Budeme přitom opět uvažovat dva případy. V prvním předpokládáme, že rozdělení pravděpodobnosti odchylek jsou pro všechny náhodné veličiny Y_1, \dots, Y_k shodná, ve druhém naopak předpokládáme, že některá jsou různá.

7.2.1. Shodná rozdělení pravděpodobnosti odchylek

Předpokládáme, že všechny odchylky ε_{ij} jsou náhodnými výběry z jediné náhodné veličiny ε . Pro testování rovnosti středních hodnot zavedeme statistiku

$$F = \frac{\sum_{i=1}^k n_i (\bar{Y}_i - \bar{Y})^2 / (k-1)}{\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2 / \left(\sum_{i=1}^k n_i - k\right)},$$

kde

$$\bar{Y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} \quad \text{a} \quad \bar{Y} = \frac{\sum_{i=1}^k n_i \bar{Y}_i}{\sum_{i=1}^k n_i}.$$

Po dosazení $\mu_i + \varepsilon_{ij}$ za Y_{ij} a při uvažování nulové hypotézy $\mu_1 = \dots = \mu_k = 0$ dostáváme

$$F_\varepsilon = \frac{\sum_{i=1}^k n_i (\bar{\varepsilon}_i - \bar{\varepsilon})^2 / (k-1)}{\sum_{i=1}^k \sum_{j=1}^{n_i} (\varepsilon_{ij} - \bar{\varepsilon}_i)^2 / \left(\sum_{i=1}^k n_i - k\right)},$$

kde

$$\bar{\varepsilon}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} \varepsilon_{ij} \quad \text{a} \quad \bar{\varepsilon} = \frac{\sum_{i=1}^k n_i \bar{\varepsilon}_i}{\sum_{i=1}^k n_i}.$$

Pokud náhodná veličina ε má normální rozdělení pravděpodobnosti, statistika F_ε má Fisherovo-Snedecorovo rozdělení pravděpodobnosti s $k-1$ a $\sum_{i=1}^k n_i - k$ stupni volnosti a platí:

$$P \left(F_\alpha < \frac{\sum_{i=1}^k n_i (\bar{\varepsilon}_i - \bar{\varepsilon})^2 / (k-1)}{\sum_{i=1}^k \sum_{j=1}^{n_i} (\varepsilon_{ij} - \bar{\varepsilon}_i)^2 / \left(\sum_{i=1}^k n_i - k\right)} < F_{1-\alpha} \right) = 1 - 2\alpha,$$

kde F_α , $F_{1-\alpha}$ jsou α -kvantil a $(1-\alpha)$ -kvantil Fisherova-Snedecorova rozdělení pravděpodobnosti s $k-1$ a $\sum_{i=1}^k n_i - k$ stupni volnosti.

Pokud nemůžeme u ε předpokládat normální rozdělení, musíme hodnoty kvantilů rozdělení pravděpodobnosti statistiky F odhadnout metodou bootstrap. Postup je následující:

1. Z pozorovaných hodnot náhodných výběrů z náhodných veličin Y_1, \dots, Y_k vypočítáme pozorované hodnoty výběrových průměrů $\bar{Y}_1, \dots, \bar{Y}_k$ a odchylek $e_{ij} = Y_{ij} - \bar{Y}_i$.

2. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu $\sum_{i=1}^k n_i$ z pozorovaných hodnot odchylek e_{11}, \dots, e_{kn_k} a označíme $e_{b,11}, \dots, e_{b,kn_k}$. Obvykle volíme $B \geq 1000$.
3. Pro každý bootstrapový výběr vypočítáme pozorované hodnoty výběrových průměrů $\bar{e}_{b,i}$ a \bar{e}_b a hodnotu statistiky $F_{e,b}$

$$F_{e,b,l} = \frac{\sum_{i=1}^k n_i (\bar{e}_{b,i} - \bar{e}_b)^2 / (k-1)}{\sum_{i=1}^k \sum_{j=1}^{n_i} (e_{b,ij} - \bar{e}_{b,i})^2 / \left(\sum_{i=1}^k n_i - k \right)},$$

kde $l = 1, \dots, B$.

4. α -kvantil a $(1 - \alpha)$ -kvantil rozdělení pravděpodobnosti statistiky $F_{e,b}$ odhadneme hodnotami $F_{b,\alpha}$ a $F_{b,1-\alpha}$ splňujícími co nejpřesněji

$$|\{F_{e,b,l}; F_{e,b,l} \leq F_{b,\alpha}\}| / B \doteq \alpha, \quad |\{F_{e,b,l}; F_{e,b,l} \leq F_{b,1-\alpha}\}| / B \doteq 1 - \alpha,$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

Pokud tedy pozorovaná hodnota statistiky F bude ležet v intervalu

$$(F_\alpha, F_{1-\alpha}), \quad \text{resp.} \quad (F_{b,\alpha}, F_{b,1-\alpha}),$$

na hladině významnosti 2α nezamítneme hypotézu $H: \mu_1 = \dots = \mu_k = 0$ a zamítneme alternativní hypotézu \bar{H} . Pokud pozorovaná hodnota statistiky F bude ležet mimo tento interval, hypotézu H naopak zamítneme a nezamítneme alternativní hypotézu \bar{H} .

7.2.2. Různá rozdělení pravděpodobnosti odchylek

V případě různých rozdělení pravděpodobnosti odchylek lze uvažovat více možností, jak vhodně upravit statistiku F . Někteří autoři uvádějí různé modifikace statistiky nebo stupňů volnosti jejího rozdělení. Mnoho příkladů z praxe ale ukazuje, že je možné i zde použít statistiku F s Fisherovým-Snedecorovým rozdělením pravděpodobnosti, jako by odchylky měly stejné normální rozdělení pravděpodobnosti a náhodné veličiny Y_1, \dots, Y_k měly stejné rozptyly. To je možné díky tomu, že statistika F je robustní a její rozdělení pravděpodobnosti je jen málo ovlivněno mírným porušením předpokladů o normalitě rozdělení odchylek nebo o shodnosti rozptylů náhodných veličin Y_1, \dots, Y_k .

Vhodným kompromisem se zdá být použití statistiky F v nezměněném tvaru, avšak s použitím metody bootstrap k odhadu hodnot kvantilů jejího rozdělení pravděpodobnosti. Postupujeme tedy takto:

1. Z pozorovaných hodnot náhodných výběrů z náhodných veličin Y_1, \dots, Y_k vypočítáme pozorované hodnoty výběrových průměrů $\bar{Y}_1, \dots, \bar{Y}_k$ a odchylek $e_{ij} = Y_{ij} - \bar{Y}_i$.
2. Realizujeme B náhodných bootstrapových výběrů (s opakováním) o rozsahu n_i z pozorovaných hodnot odchylek e_{i1}, \dots, e_{in_i} pro $i = 1, \dots, k$ a označíme je $e_{b,i1}, \dots, e_{b,in_i}$. Obvykle volíme $B \geq 1000$.

3. Pro každý bootstrapový výběr vypočítáme pozorované hodnoty výběrových průměrů $\bar{e}_{b,i}$ a \bar{e}_b a hodnotu statistiky $F_{e,b}$

$$F_{e,b,l} = \frac{\sum_{i=1}^k n_i (\bar{e}_{b,i} - \bar{e}_b)^2 / (k-1)}{\sum_{i=1}^k \sum_{j=1}^{n_i} (e_{b,ij} - \bar{e}_{b,i})^2 / \left(\sum_{i=1}^k n_i - k \right)},$$

kde $l = 1, \dots, B$.

4. α -kvantil a $(1 - \alpha)$ -kvantil rozdělení pravděpodobnosti statistiky $F_{e,b}$ odhadneme hodnotami $F_{b,\alpha}$ a $F_{b,1-\alpha}$ splňujícími co nejpřesněji

$$|\{F_{e,b,l}; F_{e,b,l} \leq F_{b,\alpha}\}| / B \doteq \alpha, \quad |\{F_{e,b,l}; F_{e,b,l} \leq F_{b,1-\alpha}\}| / B \doteq 1 - \alpha,$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

Pokud tedy pozorovaná hodnota statistiky F bude ležet v intervalu

$$(F_{b,\alpha}, F_{b,1-\alpha}),$$

na hladině významnosti 2α nezamítneme hypotézu $H: \mu_1 = \dots = \mu_k = 0$ a zamítneme alternativní hypotézu \bar{H} . Pokud pozorovaná hodnota statistiky F bude ležet mimo tento interval, hypotézu H naopak zamítneme a nezamítneme alternativní hypotézu \bar{H} .

Kapitola 8

Testy hypotéz o parametrech lineárního regresního modelu

Uvažujme lineární regresní model tak, jak byl zaveden v první kapitole, tj. $Y = X\beta + \epsilon$ a bodovým odhadem vektoru parametrů β je $b = (X^T X)^{-1} X^T Y$.

8.1. Konfidenční interval pro β_j

Konfidenční interval pro jeden libovolný koeficient β_j sestrojíme pomocí statistiky

$$t = \frac{b_j - \beta_j}{s \sqrt{((X^T X)^{-1})^{jj}}},$$

kde $((X^T X)^{-1})^{jj}$ je j -tý diagonální prvek matice $(X^T X)^{-1}$ a s je bodový odhad směrodatné odchylky podle (2.1). Pokud odchylky ϵ mají normální rozdělení pravděpodobnosti, má t Studentovo rozdělení pravděpodobnosti s $n - k$ stupni volnosti a platí

$$P \left(-t_{1-\alpha} \leq \frac{b_j - \beta_j}{s \sqrt{((X^T X)^{-1})^{jj}}} \leq t_{1-\alpha} \right) = 1 - 2\alpha,$$

odkud dostáváme konfidenční interval pro β_j se spolehlivostí $1 - 2\alpha$

$$\beta_j \in \left\langle b_j - t_{1-\alpha} s \sqrt{((X^T X)^{-1})^{jj}}; b_j + t_{1-\alpha} s \sqrt{((X^T X)^{-1})^{jj}} \right\rangle,$$

kde $t_{1-\alpha}$ je $(1 - \alpha)$ -kvantil Studentova rozdělení pravděpodobnosti s $n - k$ stupni volnosti.

Pokud ale u odchylek ϵ nemůžeme normální rozdělení předpokládat, musíme hodnoty kvantilů rozdělení pravděpodobnosti statistiky t odhadnout metodou bootstrap. Vícerozměrný bootstrapový výběr přitom budeme realizovat jako výběr z odchylek.

1. Z pozorovaných hodnot pro X a Y vypočítáme metodou nejmenších čtverců odhad b a pozorované hodnoty odchylek $e = Y - Xb$.
2. Realizujeme B bootstrapových výběrů (s opakováním) o rozsahu n z e_1, \dots, e_n a označíme $(e_{b,1,i}, \dots, e_{b,n,i})$, $i = 1, \dots, B$. Obvykle volíme $B \geq 1000$.

3. Pro každé $i = 1, \dots, B$ vypočítáme $Y_{b,i} = X\mathbf{b} + \mathbf{e}_{b,i}$. Dosazením $Y_{b,i}$ za Y dále spočítáme odhad parametrů $\mathbf{b}_{b,i}$, odhad rozptylu $s_{b,i}$ a hodnotu statistiky t_b

$$t_{b,i} = \frac{b_{b,i,j} - b_j}{s_{b,i} \sqrt{((\mathbf{X}^T \mathbf{X})^{-1})^{jj}}}.$$

4. α -kvantil a $(1-\alpha)$ -kvantil rozdělení pravděpodobnosti statistiky t_b odhadneme hodnotami $t_{b,\alpha}$, $t_{b,1-\alpha}$ splňujícími co nejpřesněji

$$|\{t_{b,i}; t_{b,i} \leq t_{b,\alpha}\}|/B \doteq \alpha, \quad |\{t_{b,i}; t_{b,i} \leq t_{b,1-\alpha}\}|/B \doteq 1 - \alpha,$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

5. *Bootstrapovým konfidenčním intervalem se spolehlivostí $1 - 2\alpha$ pro j -tý regresní parametr β_j je*

$$\beta_j \in \left\langle b_j - t_{b,1-\alpha} s \sqrt{((\mathbf{X}^T \mathbf{X})^{-1})^{jj}}; b_j - t_{b,\alpha} s \sqrt{((\mathbf{X}^T \mathbf{X})^{-1})^{jj}} \right\rangle.$$

8.2. Test hypotézy $C\boldsymbol{\beta} = \mathbf{0}$

V této podkapitole se zaměříme na problém, jak lze metodou bootstrap testovat hypotézu, která může být vyjádřena ve tvaru $H: C\boldsymbol{\beta} = \mathbf{0}$, kde C je matice typu $q \times k$ s plnou hodnotí q . Např. pro hypotézu $H: \beta_1 = \dots = \beta_k = 0$ by C byla matice typu $k \times k$

$$C = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix},$$

nebo pro hypotézu $H: \beta_1 = \beta_2$ by C byla matice typu $1 \times k$

$$C = (1 \quad -1 \quad 0 \quad \dots \quad 0).$$

Testovat budeme proti alternativní hypotéze $\overline{H}: C\boldsymbol{\beta} \neq \mathbf{0}$.

Hypotézu H testujeme pomocí statistiky

$$F = \frac{(\mathbf{C}\mathbf{b})^T [\mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T]^{-1} \mathbf{C}\mathbf{b}}{q \text{MSE}},$$

kde

$$\text{MSE} = \frac{(\mathbf{Y} - \mathbf{X}\mathbf{b})^T (\mathbf{Y} - \mathbf{X}\mathbf{b})}{n - k}. \quad (8.1)$$

Pokud odchylky $\boldsymbol{\varepsilon}$ mají normální rozdělení pravděpodobnosti a platí hypotéza H , pak statistika F má Fisherovo-Snedecorovo rozdělení pravděpodobnosti s q a $n - k$ stupni volnosti. Pokud tedy pozorovaná hodnota statistiky F bude prvkem intervalu $(F_\alpha, F_{1-\alpha})$, na hladině významnosti 2α nezamítneme hypotézu H a zamítneme alternativní hypotézu \overline{H} . Pokud pozorovaná hodnota F nebude prvkem intervalu, pak naopak zamítneme hypotézu H a nezamítneme alternativní hypotézu \overline{H} .

V případě, kdy C je řádkový vektor, tj. $q = 1$, můžeme také pro test hypotézy $H: C\beta = \Delta$ proti alternativní hypotéze $\bar{H}: C\beta \neq \Delta$ použít statistiku

$$t = \frac{Cb - C\beta}{\sqrt{\text{MSE } C(X^T X)^{-1} C^T}}.$$

Pokud odchylky $\boldsymbol{\varepsilon}$ mají normální rozdělení pravděpodobnosti, pak t má Studentovo rozdělení pravděpodobnosti s $n - k$ stupni volnosti a platí

$$P \left(-t_{1-\alpha} < \frac{Cb - C\beta}{\sqrt{\text{MSE } C(X^T X)^{-1} C^T}} < t_{1-\alpha} \right) = 1 - 2\alpha,$$

kde $t_{1-\alpha}$ je $(1 - \alpha)$ -kvantil Studentova rozdělení pravděpodobnosti s $n - k$ stupni volnosti. Odtud dostáváme konfidenční interval se spolehlivostí $1 - 2\alpha$ pro $C\beta$

$$C\beta \in \left(Cb - t_{1-\alpha} \sqrt{\text{MSE } C(X^T X)^{-1} C^T}; Cb + t_{1-\alpha} \sqrt{\text{MSE } C(X^T X)^{-1} C^T} \right).$$

Pokud Δ bude prvkem tohoto konfidenčního intervalu, na hladině významnosti 2α nezamítneme hypotézu H a zamítneme alternativní hypotézu \bar{H} . Pokud Δ nebude prvkem tohoto intervalu, pak naopak zamítneme hypotézu H a nezamítneme alternativní hypotézu \bar{H} . Tato statistika se používá zejména k testování hypotéz typu $H: \beta_j = \beta_{j0}$. Můžeme pochopitelně testovat i proti jednostranné alternativní hypotéze. V tom případě zvolíme odpovídající kvantily Studentova rozdělení.

Pokud odchylky $\boldsymbol{\varepsilon}$ nemají normální rozdělení, musíme najít vhodný tvar statistiky F , pro který můžeme odhadnout hodnoty kvantilů jejího rozdělení pravděpodobnosti metodou bootstrap. Dosadíme $\boldsymbol{\varepsilon}$ za \mathbf{Y} a označme F_ε statistiku získanou takto ze statistiky F , dále \mathbf{b}_ε odhad $\boldsymbol{\beta}$ metodou nejmenších čtverců a MSE_ε střední kvadratickou chybu podle vzorce (8.1). Pak platí

$$\mathbf{b} - \boldsymbol{\beta} = (X^T X)^{-1} X^T \mathbf{Y} - \boldsymbol{\beta} = (X^T X)^{-1} X^T (X\boldsymbol{\beta} + \boldsymbol{\varepsilon}) - \boldsymbol{\beta} = (X^T X)^{-1} X^T \boldsymbol{\varepsilon} = \mathbf{b}_\varepsilon,$$

a

$$\begin{aligned} (n - k)\text{MSE} &= (\mathbf{Y} - X\mathbf{b})^T (\mathbf{Y} - X\mathbf{b}) = (X\boldsymbol{\beta} + \boldsymbol{\varepsilon} - X\mathbf{b})^T (X\boldsymbol{\beta} + \boldsymbol{\varepsilon} - X\mathbf{b}) = \\ &= (\boldsymbol{\varepsilon} - X\mathbf{b}_\varepsilon)^T (\boldsymbol{\varepsilon} - X\mathbf{b}_\varepsilon) = (n - k)\text{MSE}_\varepsilon. \end{aligned}$$

S užitím tohoto poznatku a za předpokladu platnosti hypotézy $H: C\boldsymbol{\beta} = \mathbf{0}$ dostáváme

$$\begin{aligned} F &= \frac{(Cb)^T [C(X^T X)^{-1} C^T]^{-1} Cb}{q\text{MSE}} = \frac{[C(\mathbf{b} - \boldsymbol{\beta})]^T [C(X^T X)^{-1} C^T]^{-1} C(\mathbf{b} - \boldsymbol{\beta})}{q\text{MSE}} = \\ &= \frac{(Cb_\varepsilon)^T [C(X^T X)^{-1} C^T]^{-1} Cb_\varepsilon}{q\text{MSE}_\varepsilon} = F_\varepsilon. \end{aligned}$$

Statistiky F a F_ε tedy mají stejné rozdělení pravděpodobnosti, čehož využijeme v následujícím postupu pro test hypotézy H metodou bootstrap. Vícerozměrný bootstrapový výběr přitom opět realizujeme výběrem z odchylek.

1. Z pozorovaných hodnot pro \mathbf{X} a \mathbf{Y} vypočítáme metodou nejmenších čtverců odhad \mathbf{b} a pozorované hodnoty odchylek $\mathbf{e} = \mathbf{Y} - \mathbf{X}\mathbf{b}$.
2. Realizujeme B bootstrapových výběrů (s opakováním) o rozsahu n z e_1, \dots, e_n a označíme $(e_{b,1,i}, \dots, e_{b,n,i})$, $i = 1, \dots, B$. Obvykle volíme $B \geq 1000$.
3. Pro každé $i = 1, \dots, B$ dosazením $\mathbf{e}_{b,i}$ za \mathbf{Y} spočítáme odhad parametrů $\mathbf{b}_{b,e,i}$, střední kvadratickou chybu $\text{MSE}_{b,e,i}$ a hodnotu statistiky $F_{b,e}$

$$F_{b,e,i} = \frac{(\mathbf{C}\mathbf{b}_{b,e,i})^T [\mathbf{C}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{C}^T]^{-1}\mathbf{C}\mathbf{b}_{b,e,i}}{q\text{MSE}_{b,e,i}}$$

nebo v případě $q = 1$ hodnotu statistiky $t_{b,e}$

$$t_{b,e,i} = \frac{\mathbf{C}\mathbf{b}_{b,e,i}}{\sqrt{\text{MSE}_{b,e,i} \mathbf{C}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{C}^T}}.$$

4. α -kvantil a $(1-\alpha)$ -kvantil rozdělení pravděpodobnosti statistiky $F_{b,e}$ nebo $t_{b,e}$ odhadneme hodnotami $F_{b,\alpha}$, $F_{b,1-\alpha}$ nebo $t_{b,\alpha}$, $t_{b,1-\alpha}$ splňujícími co nejpřesněji

$$\begin{aligned} |\{F_{b,e,i}; F_{b,e,i} \leq F_{b,\alpha}\}|/B &\doteq \alpha, & |\{F_{b,e,i}; F_{b,e,i} \leq F_{b,1-\alpha}\}|/B &\doteq 1 - \alpha, \\ |\{t_{b,e,i}; t_{b,e,i} \leq t_{b,\alpha}\}|/B &\doteq \alpha, & |\{t_{b,e,i}; t_{b,e,i} \leq t_{b,1-\alpha}\}|/B &\doteq 1 - \alpha, \end{aligned}$$

kde $|\{\dots\}|$ značí velikost množiny, tj. v tomto případě počet jejích prvků.

5. Pokud pozorovaná hodnota statistiky F bude ležet v intervalu $(F_{b,\alpha}, F_{b,1-\alpha})$, na hladině významnosti 2α nezamítáme hypotézu $H: \mathbf{C}\boldsymbol{\beta} = \mathbf{0}$ a zamítáme alternativní hypotézu $\bar{H}: \mathbf{C}\boldsymbol{\beta} \neq \mathbf{0}$. Pokud pozorovaná hodnota statistiky F nebude ležet v intervalu $(F_{b,\alpha}, F_{b,1-\alpha})$, pak naopak zamítáme hypotézu H a nezamítáme alternativní hypotézu \bar{H} .
6. V případě $q = 1$ je konfidenční interval se spolehlivostí $1 - 2\alpha$ pro $\mathbf{C}\boldsymbol{\beta}$

$$\mathbf{C}\boldsymbol{\beta} \in \left(\mathbf{C}\mathbf{b} - t_{b,1-\alpha} \sqrt{\text{MSE} \mathbf{C}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{C}^T}; \mathbf{C}\mathbf{b} - t_{b,\alpha} \sqrt{\text{MSE} \mathbf{C}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{C}^T} \right).$$

Pokud Δ leží v tomto konfidenčním intervalu, nezamítáme na hladině významnosti 2α hypotézu $H: \mathbf{C}\boldsymbol{\beta} = \Delta$ a zamítáme alternativní hypotézu $\bar{H}: \mathbf{C}\boldsymbol{\beta} \neq \Delta$. Pokud Δ neleží v konfidenčním intervalu, zamítáme naopak hypotézu H a nezamítáme alternativní hypotézu \bar{H} .

8.3. Další metody

Pro testy hypotéz o regresních parametrech nebo konstruování konfidenčních intervalů pro regresní parametry můžeme využít také vícerozměrné bootstrapové výběry z $(k+1)$ -tic a kvantilové konfidenční intervaly. Doporučuje se kombinovat pivotové konfidenční intervaly s vícerozměrnými bootstrapovými výběry z odchylek (jak jsme ukázali výše v této kapitole) a kvantilové konfidenční intervaly s bootstrapovými výběry z $(k+1)$ -tic. Pak postupujeme takto:

1. Realizujeme B bootstrapových výběrů z $(k+1)$ -tic pozorovaných hodnot náhodných veličin (X_1, \dots, X_k, Y) .

2. Pro každý takový bootstrapový výběr spočítáme odhady regresních parametrů nebo jejich funkcí, které nás zajímají.
3. Odhadneme kvantily rozdělení pravděpodobnosti těchto odhadů regresních parametrů nebo jejich funkcí.
4. Sestrojíme některý z typů kvantilových konfidenčních intervalů, nejlépe BCA konfidenční interval.
5. Pomocí tohoto konfidenčního intervalu můžeme dále testovat nějakou hypotézu o regresních parametrech nebo jejich funkcích proti jednostranné nebo oboustranné alternativní hypotéze.

Kapitola 9

Odhad diskrétního rozdělení pravděpodobnosti kategoriální veličiny pomocí gradientu kvazinormy

Nechť X je kategoriální veličina, která nabývá náhodně konečně mnoha různých slovních hodnot x_j^* , $j = 1, \dots, m$, kde $m \geq 2$. Pozorováním veličiny X získáme statistický soubor (x_1, \dots, x_n) , roztríděním získáme roztríděný statistický soubor $((x_1^*, f_1/n), \dots, (x_m^*, f_m/n))$, kde $f_j/n \neq 0$ je relativní četnost pozorované hodnoty x_j^* , $j = 1, \dots, m$. Neznámé rozdělení pravděpodobnosti kategoriální veličiny X označíme $\mathbf{p} = (p_1, \dots, p_m)$, kde $p_j = P(X = x_j^*)$. Odhad pravděpodobností \mathbf{p} je vlastně odhadem parametrů multinomického rozdělení pravděpodobnosti $M(n, p_1, \dots, p_m)$. Nestranným odhadem \mathbf{p} je $\hat{\mathbf{p}} = (f_1/n, \dots, f_m/n)$. V [7] je předložen následující pesimistický odhad \mathbf{p} založený na gradientu kvazinormy rozdělení \mathbf{p} .

Nechť funkce $f: (0, \infty) \rightarrow \mathbb{R}^*$, kde $\mathbb{R}^* = \mathbb{R} \cup \{-\infty, \infty\}$, je konvexní na $(0, \infty)$, striktně konvexní v bodě $u = 1$ a nabývá zde hodnoty $f(1) = 0$. Jestliže $\mathbf{p} = (p_1, \dots, p_m)$, resp. $\mathbf{q} = (q_1, \dots, q_m)$ je diskrétní rozdělení pravděpodobnosti z pravděpodobnostního prostoru (Ω, Σ, P) , resp. (Ω, Σ, Q) , pak f -divergencí těchto rozdělení rozumíme funkcionál

$$D_f(\mathbf{p}, \mathbf{q}) = \sum_{j=1}^m q_j f\left(\frac{p_j}{q_j}\right).$$

Pojem f -divergence má význam vzdálenosti daných rozdělení. Platí

1. $\mathbf{p} = \mathbf{q} \Leftrightarrow D_f(\mathbf{p}, \mathbf{q}) = 0$,
2. $D_f(\mathbf{p}, \mathbf{q})$ nabývá v \mathbb{R}^* svého maxima $\Leftrightarrow \mathbf{p}$ a \mathbf{q} jsou ortogonální, tj. existují takové disjunktní množiny $E, F \subset \Omega$, že

$$\sum_{x_j^* \in E} p_j = 1 \quad \text{a} \quad \sum_{x_j^* \in F} q_j = 1.$$

Nechť $S = \left\{ \mathbf{p} \in \mathbb{R}^m : \forall p_j \geq 0, \sum_{j=1}^m p_j = 1 \right\}$ je množina všech diskrétních rozdělení pravděpodobnosti na Ω . Kvazinormou rozdělení $\mathbf{p} \in S$ rozumíme f -divergenci $D_f(\mathbf{p}, \mathbf{p}_0)$, kde $\mathbf{p}_0 = (1/m, \dots, 1/m)$, a o funkci f říkáme, že generuje kvazinormu $D_f(\mathbf{p}, \mathbf{p}_0)$ na S . Platí

1. $D_f(\mathbf{p}, \mathbf{p}_0) = \frac{1}{m} \sum_{j=1}^m f(mp_j)$,
2. $D_f(\mathbf{p}, \mathbf{p}_0)$ je nezáporná konvexní funkce na S symetrická vzhledem k proměnným p_j , $j = 1, \dots, m$.
3. \mathbf{p}_0 minimalizuje integrál všech f -divergencí $D_f(\mathbf{p}, \mathbf{q})$ na S a má maximální entropii.

Budeme hledat takové rozdělení pravděpodobnosti v S , které je nejbližší \mathbf{p}_0 a k němuž se dostaneme od empirického rozdělení co nejrychleji. Tomu odpovídá minimalizace kvazinormy $D_f(\mathbf{p}, \mathbf{p}_0)$ a hledání rozdělení na křivce největšího spádu v S .

Nechť $D_f(\mathbf{p}, \mathbf{p}_0)$ je kvazinorma na S . *Gradientním odhadem* rozdělení pravděpodobnosti $\mathbf{p} \in S$ z empirického rozdělení $(f_1/n, \dots, f_m/n)$ rozumíme takové rozdělení pravděpodobnosti $\mathbf{p}(t) \in S$, že

$$\frac{d}{dt} \mathbf{p}(t) = -\text{grad } D_f(\mathbf{p}(t), \mathbf{p}_0) \quad \forall t \in \langle 0, \infty \rangle \quad \text{a} \quad \mathbf{p}(0) = \mathbf{f}/n = (f_1/n, \dots, f_m/n).$$

Jestliže funkce $f(u)$ generuje kvazinormu $D_f(\mathbf{p}, \mathbf{p}_0)$ na S , má výše uvedené vlastnosti a má spojitou derivaci $f'(u)$ pro každé $u \in (0, \infty)$, pak existuje jediný gradientní odhad $\mathbf{p}(t) = (p_1(t), \dots, p_m(t))$ rozdělení pravděpodobnosti $\mathbf{p} \in S$. Jeho složky $p_1(t), \dots, p_{m-1}(t)$ jsou $\forall t \in \langle 0, \infty \rangle$ partikulárním řešením soustavy obyčejných diferenciálních rovnic prvního řádu

$$\begin{aligned} p_1'(t) &= -f'(mp_1(t)) + f' \left(m \left[1 - \sum_{j=1}^{m-1} p_j(t) \right] \right), \\ &\vdots \\ p_{m-1}'(t) &= -f'(mp_{m-1}(t)) + f' \left(m \left[1 - \sum_{j=1}^{m-1} p_j(t) \right] \right) \end{aligned}$$

s počátečními podmínkami

$$p_1(0) = f_1/n, \dots, p_{m-1}(0) = f_{m-1}/n$$

a složka $p_m(t) = 1 - \sum_{j=1}^{m-1} p_j(t) \quad \forall t \in \langle 0, \infty \rangle$.

Vhodnou hodnotu parametru $t_0 \in \langle 0, \infty \rangle$ nalezneme pomocí testu dobré shody jako takovou hodnotu t , kdy ještě nezamítáme hypotézu o vhodnosti rozdělení $\mathbf{p}(t)$ na hladině významnosti α . Gradientní odhad $\mathbf{p}(t)$ se pro rostoucí parametr t vzdaluje po křivce největšího spádu v S od empirického rozdělení k \mathbf{p}_0 . Odhad $\mathbf{p}(t_0)$ je nejhorším z odhadů, které splňují zvolené testové kritérium na hladině významnosti alespoň α , proto ho nazýváme *pesimistickým gradientním odhadem*.

Nechť $f(u) = (u - 1)^2$. Pak $D_f(\mathbf{p}, \mathbf{p}_0) = \frac{1}{m} \sum_{j=1}^m (mp_j - 1)^2$ je tzv. *kvadratická kvazinorma*. Potom složky gradientního odhadu $\mathbf{p}(t) = (p_1(t), \dots, p_m(t))$ z empirického rozdělení $(f_1/n, \dots, f_m/n)$ jsou $\forall t \in \langle 0, \infty \rangle$ partikulárním řešením nehomogenní lineární soustavy obyčejných diferenciálních rovnic prvního řádu s konstantními koeficienty a pravými stranami

$$\begin{aligned} p_1'(t) &= -4mp_1(t) - 2mp_2(t) - \dots - 2mp_{m-1}(t) + 2m, \\ p_2'(t) &= -2mp_1(t) - 4mp_2(t) - \dots - 2mp_{m-1}(t) + 2m, \\ &\vdots \\ p_{m-1}'(t) &= -2mp_1(t) - 2mp_2(t) - \dots - 4mp_{m-1}(t) + 2m \end{aligned}$$

s počátečními podmínkami

$$p_1(0) = f_1/n, \dots, p_{m-1}(0) = f_{m-1}/n$$

a složka $p_m(t) = 1 - \sum_{j=1}^{m-1} p_j(t) \quad \forall t \in \langle 0, \infty \rangle$.

Řešením soustavy získáváme složky gradientního odhadu $\mathbf{p}(t)$

$$\begin{aligned} p_1(t) &= c_1 e^{-2m^2 t} + c_2 e^{-2mt} && + 1/m, \\ p_2(t) &= c_1 e^{-2m^2 t} && + c_3 e^{-2mt} + 1/m, \\ &\vdots \\ p_{m-2}(t) &= c_1 e^{-2m^2 t} && + c_{m-1} e^{-2mt} + 1/m, \\ p_{m-1}(t) &= c_1 e^{-2m^2 t} - c_2 e^{-2mt} \dots - c_{m-1} e^{-2mt} && + 1/m, \\ p_m(t) &= -(m-1) c_1 e^{-2m^2 t} && + 1/m, \end{aligned}$$

kde

$$\begin{aligned} c_1 &= \frac{f_1/n + f_2/n + \dots + f_{m-1}/n}{m-1} - \frac{1}{m}, \\ c_2 &= \frac{(m-2)f_1/n - f_2/n - \dots - f_{m-1}/n}{m-1}, \\ c_3 &= \frac{-f_1/n + (m-2)f_2/n - \dots - f_{m-1}/n}{m-1}, \\ &\vdots \\ c_{m-1} &= \frac{-f_1/n - f_2/n - \dots + (m-2)f_{m-2}/n - f_{m-1}/n}{m-1}. \end{aligned}$$

Složky takto získaného gradientního odhadu z empirického rozdělení jsou asymptoticky nestrannými odhady složek pozorovaného rozdělení pravděpodobnosti \mathbf{p} .

Výpočet hodnoty parametru t_0 je dosti citlivý na „strmost“ zvolené kvazinormy, související numerické problémy s řešením odpovídajících nelineárních diferenciálních rovnic lze alespoň částečně zmenšit vhodným kladným násobkem kvazinormy. Kvadratická kvazinorma je přitom až na tento násobek jediná, která vede na lineární soustavu diferenciálních rovnic s konstantními koeficienty.

Kapitola 10

Intervalové odhady diskrétního rozdělení pravděpodobnosti kategoriální veličiny

V minulé kapitole jsme uvedli dva odhady diskrétního rozdělení pravděpodobnosti kategoriální veličiny, a to odhad relativními četnostmi a pesimistický gradientní odhad pomocí kvadratické kvazinormy (dále zkráceně jen pesimistický gradientní odhad). V obou případech se ale jedná pouze o bodové odhady. Jednou z možností, jak získat intervalové odhady, je použití metody bootstrap. V této kapitole uvedeme dva příklady kategoriálních veličin, u nichž předvedeme získání intervalových odhadů jejich rozdělení pravděpodobnosti metodou bootstrap. Dále provedeme srovnání intervalových odhadů získaných pomocí obou výše uvedených bodových odhadů a také ukážeme vliv rozsahu pozorovaného náhodného výběru, stejně jako vliv počtu bootstrapových výběrů B .

K simulacím a výpočtům uvedeným v této kapitole jsme použili matematický software Shine bootstrap a Matlab R2008a. Shine bootstrap byl vytvořen speciálně pro generování náhodných výběrů z diskrétních rozdělení pravděpodobnosti a pro výpočty bodových odhadů pravděpodobnostních funkcí těchto rozdělení. Zejména je v něm implementován výpočet pesimistického gradientního odhadu podle postupu předloženého v předchozí kapitole. Matlab R2008a byl použit zejména k výpočtu kvantilů a generování histogramů.

10.1. Falešná kostka

V prvním příkladě budeme uvažovat hrací kostku o šesti stranách, které si označíme čísly $1, \dots, 6$. Házíme-li kostkou, pozorujeme diskrétní náhodnou veličinu X — číslo, které padne. Základní prostor je tedy tvořen šesti elementárními náhodnými jevy, které odpovídají číslům $1, \dots, 6$. V případě, kdy kostka není falešná, je pravděpodobnostní funkce této náhodné veličiny $\mathbf{p} = (p_1, \dots, p_6) = (1/6, \dots, 1/6) \doteq (0,1667; \dots; 0,1667)$. My ale budeme uvažovat kostku falešnou, např. takovou, která má poněkud těžší stranu s číslem 6. Její pravděpodobnostní funkci zvolíme $\mathbf{p} = (p_1, \dots, p_6) = (0,10; 0,15; 0,15; 0,15; 0,15; 0,30)$. Pozorování náhodné veličiny X realizujeme na počítači jako náhodný výběr z diskrétního rozdělení se zadanou pravděpodobnostní funkcí. Protože chceme mimo jiné zkoumat vliv rozsahu n pozorovaného náhodného výběru, realizujeme tři pozorování a získáme tři náhodné výběry z veličiny X , a to o rozsazích 60, 100 a 400. Pro všechny tři náhodné výběry vypočítáme bodové

	strana 1	strana 2	strana 3	strana 4	strana 5	strana 6	rozsah n
skutečné p	0,10	0,15	0,15	0,15	0,15	0,30	
odhad	0,1167	0,1333	0,1333	0,0833	0,1667	0,3667	60
relativními	0,1600	0,1400	0,1800	0,1100	0,1000	0,3100	100
četnostmi	0,1200	0,1425	0,1525	0,1350	0,1500	0,3000	400
pesimistický	0,1235	0,1397	0,1397	0,0911	0,1720	0,3341	60
gradientní	0,1642	0,1448	0,1835	0,1157	0,1060	0,2859	100
odhad	0,1228	0,1449	0,1548	0,1375	0,1523	0,2877	400

Tab. 10.1: Bodové odhady pravděpodobnostní funkce falešné kostky

odhady pravděpodobnostní funkce p , a to oběma výše uvedenými způsoby, tj. bodový odhad relativními četnostmi i bodový pesimistický gradientní odhad. Výsledky jsou zaznamenány v tabulce 10.1.

Vidíme, že pro strany 1, ..., 5 je vždy bodový pesimistický gradientní odhad vyšší než příslušný bodový odhad relativními četnostmi, zatímco pro stranu 6 je vždy nižší. Pesimistický gradientní odhad má tedy tendenci snižovat rozdíly mezi jednotlivými pravděpodobnostmi oproti pozorovaným relativním četnostem. Pro strany 2, ..., 6 jsme obdrželi nejpřesnější bodové odhady jednotlivých pravděpodobností pro náhodný výběr o rozsahu 400, u všech stran 1, ..., 6 jsou odhady jednotlivých pravděpodobností pro náhodný výběr o rozsahu 400 blíže skutečným hodnotám než odhady pro náhodný výběr o rozsahu 100. Z toho lze odvodit, že přesnost bodových odhadů stoupá s rozsahem pozorovaného náhodného výběru. Srovnáním přesnosti všech 18 bodových odhadů jednotlivých pravděpodobností získaných jako odhad relativními četnostmi a jako pesimistický gradientní odhad zjišťujeme, že v 9 případech byl přesnější odhad relativními četnostmi a v 9 případech byl přesnější pesimistický gradientní odhad. Bodové odhady získané oběma metodami jsou tedy zhruba stejně přesné, žádná z metod se nejeví být lepší. Bude nás zajímat, zda bude pozorovatelný nějaký rozdíl přesnosti obou metod u intervalových odhadů.

Nejjednodušší a výpočtově nejméně náročný způsob zisku intervalového odhadu je v tomto případě jednoduchý kvantilový bootstrapový konfidenční interval, který nevyžaduje žádné dodatečné výpočty směrodatných odchylek, vychýlení mediánu, akcelerace ani reziduí. Jednoduše realizujeme B bootstrapových výběrů z pozorovaného náhodného výběru o rozsahu n a pro každý z nich vypočítáme bodový odhad pravděpodobnostní funkce. Pro každou stranu kostky (každý elementární náhodný jev náhodné veličiny X) tak dostaneme B hodnot (bodových odhadů). Jejich α -kvantil a $(1 - \alpha)$ -kvantil pak jsou krajními body intervalového odhadu se spolehlivostí $1 - 2\alpha$. My přitom celý postup zopakujeme pro všechny tři náhodné výběry o rozsazích 60, 100 a 400 a také pro oba odhady pravděpodobnostní funkce, tj. pro odhad relativními četnostmi i pro pesimistický gradientní odhad. Protože chceme dále zkoumat i vliv počtu bootstrapových výběrů, zvolíme také dvě různé hodnoty B , a to 1000 a 5000. Výsledky jsou zaznamenány v tabulce 10.2, volili jsme spolehlivost 0,95.

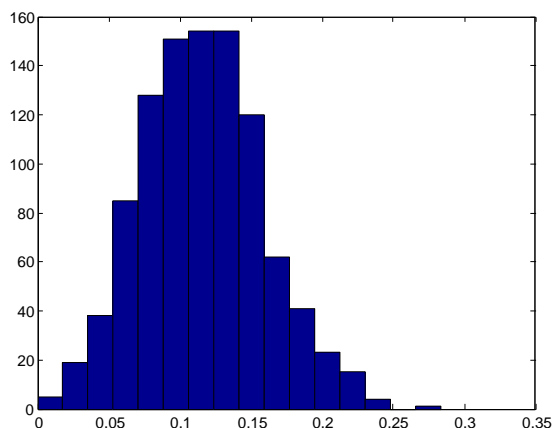
Celkem je tedy v tabulce 72 intervalových odhadů. Každý z těchto intervalových odhadů přitom obsahuje příslušnou skutečnou hodnotu pravděpodobnosti daného náhodného jevu. Dále je ve všech případech splněno, že šířka intervalu klesá se stoupajícím rozsahem pozorovaného náhodného výběru. Ve 36 případech můžeme srovnat šířku intervalu získaného pomocí

	strana 1		strana 2		strana 3		rozsah n	počet B
skutečné p	0,10		0,15		0,15			
odhad relativními četnostmi	0,0500	0,2000	0,0500	0,2167	0,0500	0,2167	60	1000
	0,0500	0,2000	0,0500	0,2167	0,0500	0,2167	60	5000
	0,0900	0,2400	0,0800	0,2100	0,1200	0,2700	100	1000
	0,0900	0,2300	0,0800	0,2100	0,1200	0,2700	100	5000
	0,0900	0,1525	0,1100	0,1800	0,1187	0,1900	400	1000
	0,0900	0,1525	0,1075	0,1775	0,1175	0,1900	400	5000
pesimistický gradientní odhad	0,0573	0,2042	0,0585	0,2205	0,0580	0,2354	60	1000
	0,0572	0,2043	0,0583	0,2206	0,0582	0,2206	60	5000
	0,0959	0,2411	0,0769	0,2128	0,1155	0,2609	100	1000
	0,0960	0,2318	0,0859	0,2127	0,1154	0,2612	100	5000
	0,0931	0,1548	0,1129	0,1792	0,1202	0,1917	400	1000
	0,0931	0,1548	0,1128	0,1793	0,1203	0,1892	400	5000
	strana 4		strana 5		strana 6		rozsah n	počet B
skutečné p	0,15		0,15		0,30			
odhad relativními četnostmi	0,0167	0,1667	0,0833	0,2667	0,2500	0,4833	60	1000
	0,0167	0,1667	0,0833	0,2667	0,2500	0,4833	60	5000
	0,0500	0,1700	0,0400	0,1600	0,2200	0,4100	100	1000
	0,0500	0,1700	0,0400	0,1600	0,2200	0,4000	100	5000
	0,1025	0,1725	0,1150	0,1875	0,2550	0,3425	400	1000
	0,1025	0,1700	0,1175	0,1875	0,2575	0,3450	400	5000
pesimistický gradientní odhad	0,0241	0,1558	0,0751	0,2674	0,2230	0,4670	60	1000
	0,0242	0,1559	0,0901	0,2676	0,2219	0,4662	60	5000
	0,0567	0,1830	0,0561	0,1641	0,2005	0,3803	100	1000
	0,0566	0,1800	0,0475	0,1641	0,2019	0,3750	100	5000
	0,1055	0,1719	0,1178	0,1843	0,2422	0,3347	400	1000
	0,1078	0,1719	0,1178	0,1868	0,2434	0,3323	400	5000

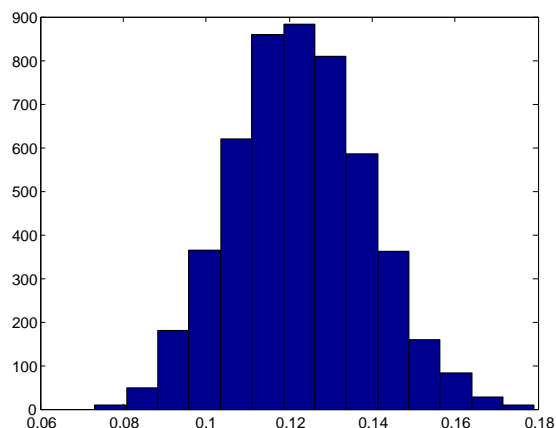
Tab. 10.2: Bootstrapové intervalové odhady pravděpodobnostní funkce falešné kostky

1000 bootstrapových výběrů a pomocí 5000 bootstrapových výběrů. Pouze v 8 případech jsme pro $B = 5000$ obdrželi širší intervalový odhad než pro $B = 1000$, přičemž největší rozšíření intervalu bylo o 0,0086. V 15 případech byly oba intervaly stejně široké, v ostatních 13 případech se interval s rostoucím B zúžil. Ve 36 případech můžeme také srovnat šířku intervalového odhadu relativními četnostmi a intervalového pesimistického gradientního odhadu. Ve 26 z nich byl intervalový pesimistický gradientní odhad užší.

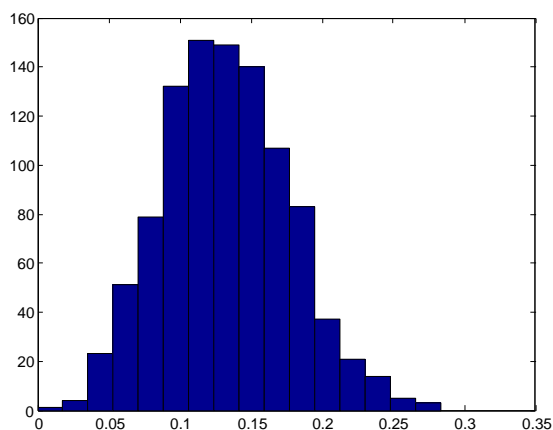
Na následujících stranách uvádíme 12 histogramů B bodových odhadů pro jednotlivé strany kostky (jednotlivé náhodné jevy). Pro každou ze stran 1, ..., 6 jsou to vždy případy $n = 60$ a $n = 400$, kterým odpovídají nejširší a nejužší ze všech intervalových odhadů pro tu kterou stranu kostky. U každého histogramu je uvedena hodnota B , a zda se v daném případě jedná o intervalový odhad relativními četnostmi nebo intervalový pesimistický gradientní odhad.



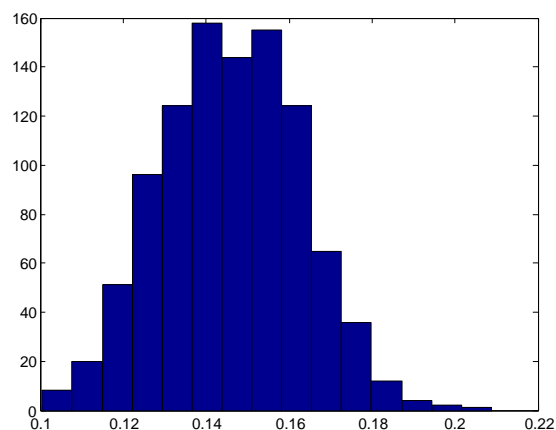
Strana 1, $n = 60$, $B = 1000$,
odhad relativními četnostmi



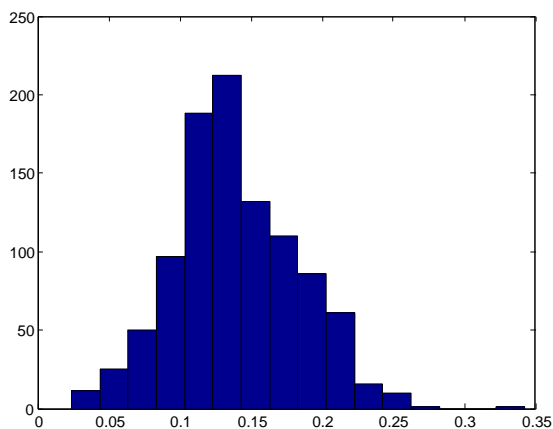
Strana 1, $n = 400$, $B = 5000$,
pesimistický gradientní odhad



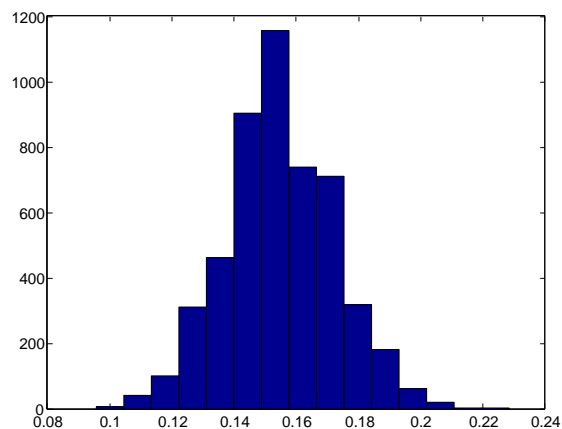
Strana 2, $n = 60$, $B = 1000$,
odhad relativními četnostmi



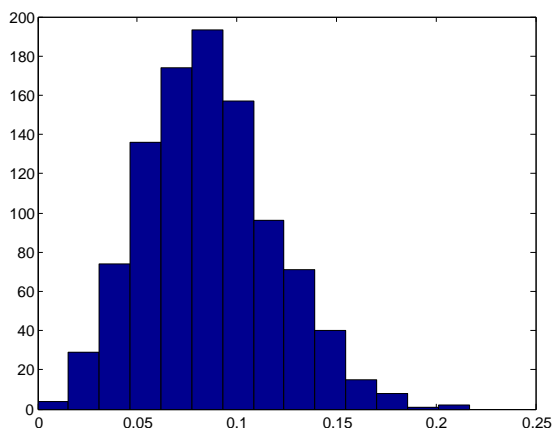
Strana 2, $n = 400$, $B = 1000$,
pesimistický gradientní odhad



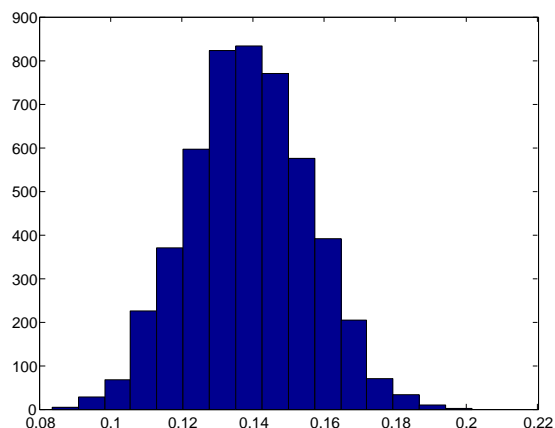
Strana 3, $n = 60$, $B = 1000$,
pesimistický gradientní odhad



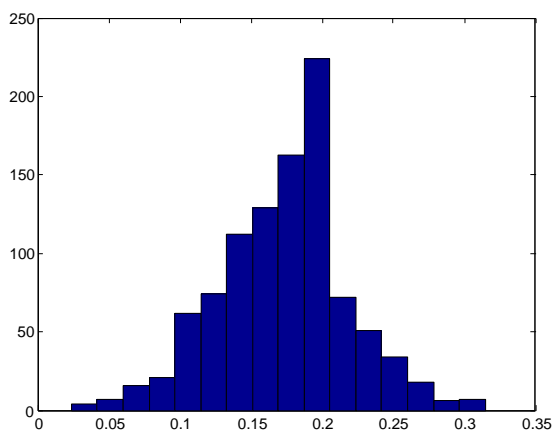
Strana 3, $n = 400$, $B = 5000$,
pesimistický gradientní odhad



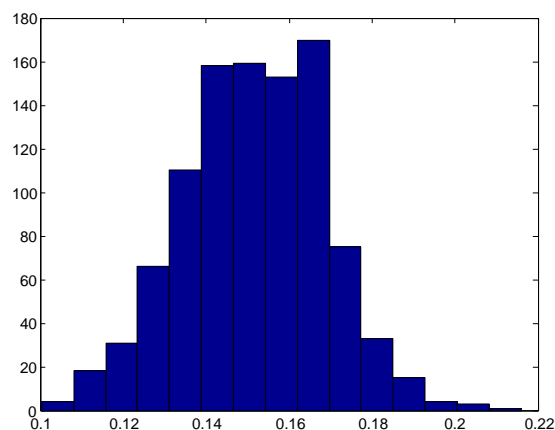
Strana 4, $n = 60$, $B = 1000$,
odhad relativními četnostmi



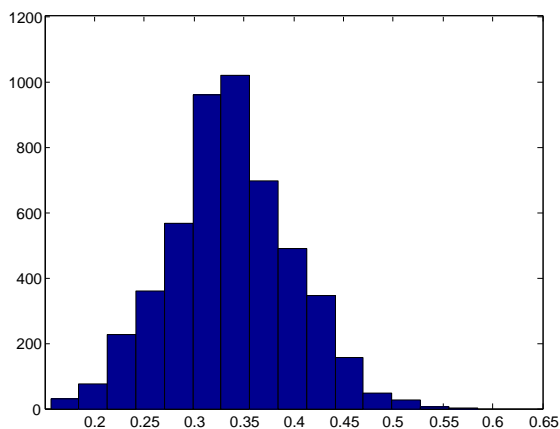
Strana 4, $n = 400$, $B = 5000$,
pesimistický gradientní odhad



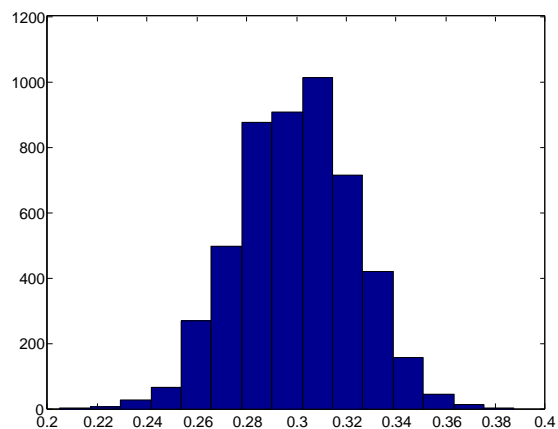
Strana 5, $n = 60$, $B = 1000$,
pesimistický gradientní odhad



Strana 5, $n = 400$, $B = 1000$,
pesimistický gradientní odhad



Strana 6, $n = 60$, $B = 5000$,
pesimistický gradientní odhad



Strana 6, $n = 400$, $B = 5000$,
odhad relativními četnostmi

Podstatnou otázkou také je, zda jsou získané intervalové odhady pravděpodobnostní funkce dostatečně úzké, abychom na základě nich byli schopni rozhodnout, zda kostka je nebo není falešná. Testujme hypotézu $H: p_6 = 1/6$ proti alternativní hypotéze $\bar{H}: p_6 \neq 1/6$. Z tabulky je patrné, že již pro náhodný výběr o rozsahu 60 (a dále i pro oba větší rozsahy) žádný ze čtyř intervalových odhadů pro stranu s číslem 6 neobsahuje hodnotu 0,1667. Tedy na hladině významnosti 0,05 zamítáme hypotézu H a nezamítáme alternativní hypotézu \bar{H} . Z toho ovšem vyplývá, že kostka je falešná. Vzhledem k tomu, že pro $\alpha < 0,5$ platí, že 0-kvantil $\leq \alpha$ -kvantil $\leq 2\alpha$ -kvantil $\leq (1 - 2\alpha)$ -kvantil $\leq (1 - \alpha)$ -kvantil ≤ 1 -kvantil, můžeme bez dalších výpočtů říci, že také při testu hypotézy H proti alternativní hypotéze $\bar{H}: p_6 \geq 1/6$ zamítáme H a nezamítáme \bar{H} . Strana s číslem 6 padá na této kostce příliš často.

Dále pro náhodný výběr o rozsahu 100 vidíme, že ani žádný ze čtyř intervalových odhadů pro stranu s číslem 5 neobsahuje hodnotu 0,1667, pro náhodný výběr o rozsahu 400 žádný ze čtyř intervalových odhadů pro stranu s číslem 1 neobsahuje hodnotu 0,1667. Podobně jako v předchozím odstavci dojdeme k závěru, že strany s čísly 1 a 5 naopak padají příliš málo.

Závěrem lze říci, že metodu bootstrap lze pro získání intervalových odhadů pravděpodobnostní funkce této kategoriální veličiny doporučit, získané intervalové odhady jsou dostatečně přesné a úzké. Dalšího zlepšení lze dosáhnout větším rozsahem pozorovaného náhodného výběru, větším počtem bootstrapových výběrů a použitím metody pesimistického gradientního odhadu.

10.2. Volební model

Ve druhém příkladu se budeme věnovat předvolebnímu průzkumu veřejného mínění. Ve dnech 17. až 18. října 2008 se konaly volby do krajských zastupitelstev. Hlasy voličů se v tomto případě sčítají pro každý kraj samostatně, zaměříme se proto např. na Jihomoravský kraj. Ve dnech 5. až 7. října 2008 realizovala společnost CS&C pro Českou televizi tzv. bleskový předvolební průzkum veřejného mínění (viz [10]) občanů Jihomoravského kraje. Respondentů bylo 633, z nich pouze 306 bylo zahrnuto do tzv. volebního modelu. Do volebního modelu se zahrnují pouze odpovědi těch respondentů, kteří chtějí jít volit a zároveň už vědí, koho budou volit. K volbám se nakonec v Jihomoravském kraji dostavilo 377 706 voličů (jejichž hlasy byly platné — viz [9]). Kandidovalo celkem 18 politických stran a hnutí. Ve volebním modelu je uvedeno 8 z nich, které měly nejvyšší preference, všechny ostatní tvoří devátou položku modelu.

Volební model je tedy kategoriální veličina X , která může nabývat devíti slovních hodnot (elementárních náhodných jevů), a to ODS, ČSSD, KSČM, KDU-ČSL, SZ, SNK-ED, NEZÁVISLÍ, Moravané a ostatní. Pozorováním kategoriální veličiny X byl získán náhodný výběr o rozsahu 306. V závěrečné zprávě průzkumu jsou uvedeny absolutní četnosti pozorovaných hodnot a na jejich základě vytvořený odhad procentuálního zisku jednotlivých politických stran v nadcházejících volbách. Je uveden bodový odhad relativními četnostmi a dále i intervalový odhad se spolehlivostí 0,95, není ale zveřejněno, jakou metodou byl získán. K jeho výpočtu nicméně byla použita asymptotická metoda, která předpokládá absolutní četnost pozorování každého elementárního náhodného jevu alespoň 5. Tento předpoklad ale nebyl splněn.

My odhadneme složky pravděpodobnostní funkce p kategoriální veličiny X oběma výše uvedenými bodovými odhady, tj. realizujeme odhad relativními četnostmi a pesimistický gradientní odhad. Pro získání intervalových odhadů použijeme opět jednoduchý kvantilový bootstrap

pový konfidenční interval. Realizujeme B bootstrapových výběrů z pozorovaného náhodného výběru a pro každý z nich vypočítáme odhad pravděpodobnostní funkce. Pro každý elementární náhodný jev tak dostaneme B hodnot (bodových odhadů). Jejich α -kvantil a $(1 - \alpha)$ -kvantil pak jsou krajními body intervalového odhadu se spolehlivostí $1 - 2\alpha$. Celý postup přitom zopakujeme pro oba bodové odhady, tj. odhad relativními četnostmi i pesimistický gradientní odhad. Zajímá nás také vliv počtu bootstrapových výběrů, proto budeme volit dvě různé hodnoty B , a to 1000 a 5000. Spolehlivost volíme 0,95, abychom mohli provést relevantní srovnání s intervalovými odhady uvedenými v [10]. Pozorovaný náhodný výběr, odhady převzaté z [10], výsledky provedených výpočtů i skutečný výsledek voleb jsou uvedeny v tabulce 10.3.

Vidíme, že pro prvních pět elementárních náhodných jevů, jejichž relativní četnosti přesahují 0,05, je bodový pesimistický gradientní odhad nižší než bodový odhad relativními četnostmi, zatímco pro následující čtyři elementární jevy, jejichž relativní četnosti nedosahují ani 0,03, je bodový pesimistický gradientní odhad vyšší než bodový odhad relativními četnostmi. Opět se tedy potvrzuje, že pesimistický gradientní odhad má tendenci snižovat rozdíly mezi jednotlivými pravděpodobnostmi oproti pozorovaným relativním četnostem.

V tabulce je zaznamenáno celkem 45 intervalových odhadů, z toho 9 jsme převzali z [10] a 36 jsme vypočítali metodou bootstrap, a to 18 pomocí odhadu relativními četnostmi a 18 pomocí pesimistického gradientního odhadu. Můžeme tedy právě v 18 případech srovnat intervalový odhad relativními četnostmi s intervalovým pesimistickým gradientním odhadem. Ve 12 případech je intervalový pesimistický gradientní odhad užší než příslušný intervalový odhad relativními četnostmi, v 1 případě byly oba intervaly stejně široké a pouze v 5 případech je intervalový pesimistický gradientní odhad širší než příslušný intervalový odhad relativními četnostmi. Opět tedy vidíme, že intervalové pesimistické gradientní odhady jsou přesnější. Dále je patrné, že pro prvních 5 elementárních náhodných jevů jsou intervalové pesimistické gradientní odhady oproti intervalovým odhadům relativními četnostmi posunuté doleva, zatímco pro zbylé 4 elementární jevy naopak doprava, čímž se opět potvrzuje, že pesimistický gradientní odhad má tendenci snižovat rozdíly mezi jednotlivými pravděpodobnostmi oproti pozorovaným relativním četnostem.

V 18 případech také můžeme srovnat intervaly získané pomocí $B = 1000$ a $B = 5000$ bootstrapových výběrů. V 8 případech se interval s rostoucím B zúžil, v 6 případech jsme získali zcela stejné (tudíž i stejně široké) intervaly a ve 4 případech se interval s rostoucím B rozšířil, a to nejvýše o 0,0049. Vyšší počet bootstrapových výběrů má tedy i v tomto příkladu kategoriální veličiny pozitivní vliv na intervalové odhady pravděpodobnostní funkce.

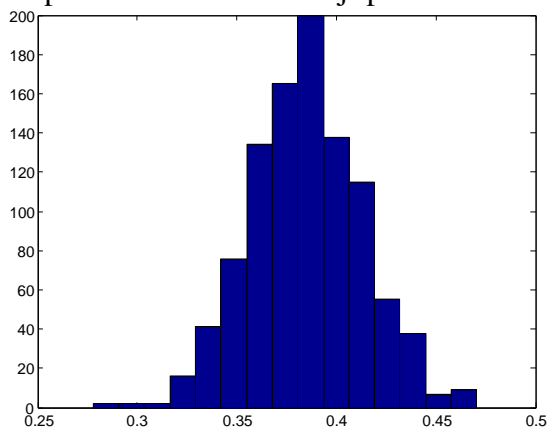
Na následujících stranách uvádíme 18 histogramů B bodových odhadů pro jednotlivé elementární náhodné jevy. Pro každý z jevů 2 histogramy, které odpovídají nejširšímu a nejužšímu intervalovému odhadu pro daný jev. U každého histogramu je uvedeno, zda se v daném případě jedná o intervalový odhad relativními četnostmi, nebo intervalový pesimistický gradientní odhad, a pomocí kolika bootstrapových výběrů byl získán.

Dále je třeba srovnat převzaté intervalové odhady s těmi, které jsme vypočítali metodou bootstrap. Krajní meze převzatých intervalových odhadů byly autory zaokrouhleny na celá procenta, tj. na dvě desetinná místa. Pro prvních šest elementárních náhodných jevů lze říci, že pokud bychom zaokrouhlili hodnoty krajních mezí bootstrapových intervalových odhadů na dvě desetinná místa, obdrželi bychom takřka vždy, konkrétně ve 45 případech z 48, odhady převzaté z [10], ve zbylých 3 případech hodnotu o 0,01 vyšší. U zbylých tří elementárních jevů

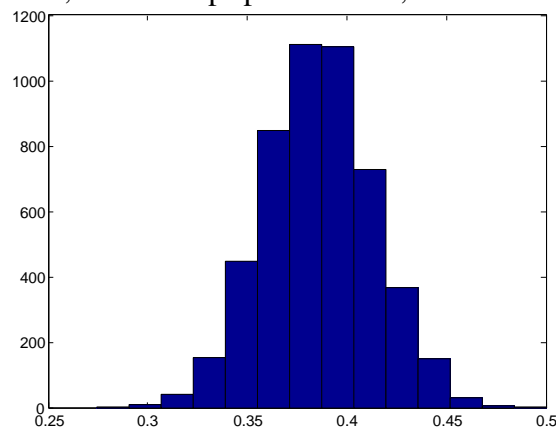
	ČSSD		ODS		KSČM		KDU-ČSL		SZ	počet <i>B</i>
absolutní četnosti	118	71	35	52	17					
odhad SC&C	0,33	0,44	0,18	0,28	0,8	0,15	0,13	0,21	0,3	0,8
odhad relativními četnostmi	0,3856	0,2320	0,1144	0,1699					0,0556	
pesimistický gradientní odhad	0,3815	0,2296	0,1133	0,1682					0,0551	
intervalový odhad	0,3301	0,4444	0,1863	0,2810	0,0784	0,1503	0,1291	0,2124	0,0294	0,0817
relativními četnostmi	0,3333	0,4412	0,1863	0,2810	0,0784	0,1503	0,1275	0,2157	0,0327	0,0817
intervalový pesimistický	0,3270	0,4389	0,1846	0,2788	0,0778	0,1490	0,1277	0,2113	0,0293	0,0812
gradientní odhad	0,3279	0,4368	0,1842	0,2776	0,0780	0,1494	0,1265	0,2128	0,0325	0,0810
výsledky voleb	0,3484	0,1588	0,1441	0,2389					0,0364	
	NEZÁVISLÍ		Moravané		SNK-ED		ostatní			počet <i>B</i>
absolutní četnosti	3	4	2	4						
odhad SC&C	0,00	0,02	0,00	0,02	0,00	0,01	0,00	0,02		
odhad relativními četnostmi	0,0098	0,0131	0,0065	0,0131						
pesimistický gradientní odhad	0,0099	0,0132	0,0067	0,0132						
intervalový odhad	0,0000	0,0229	0,0033	0,0294	0,0000	0,0163	0,0033	0,0261		1000
relativními četnostmi	0,0000	0,0229	0,0030	0,0261	0,0000	0,0163	0,0033	0,0261		5000
intervalový pesimistický	0,0002	0,0228	0,0033	0,0293	0,0001	0,0165	0,0093	0,0382		1000
gradientní odhad	0,0002	0,0228	0,0033	0,0261	0,0001	0,0164	0,0093	0,0383		5000
výsledky voleb	0,0189	0,0089	0,0139	0,0417						

Tab. 10.3: Volební model a výsledky voleb do krajských zastupitelstev 2008

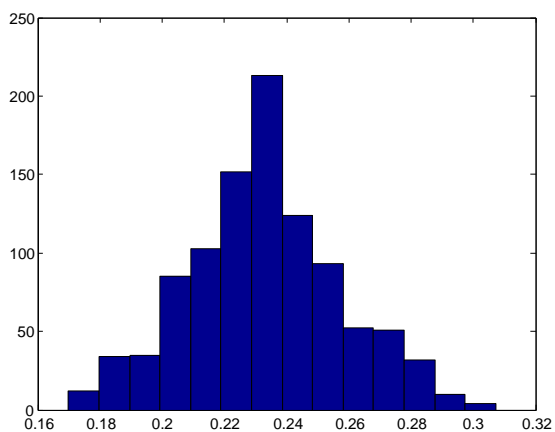
se po zaokrouhlení shodují pouze dolní meze intervalů, a to v 10 případech z 12, horní meze



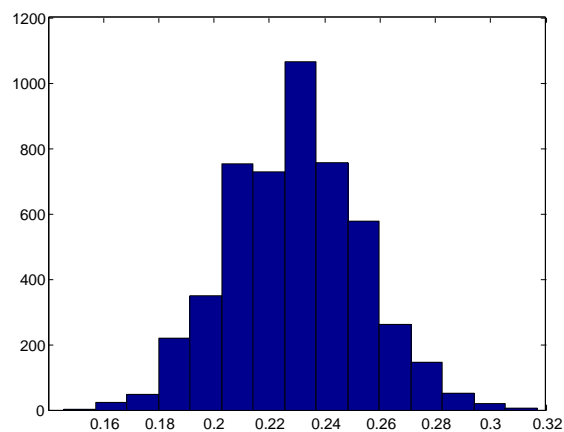
ČSSD, $B = 1000$,
odhad relativními četnostmi



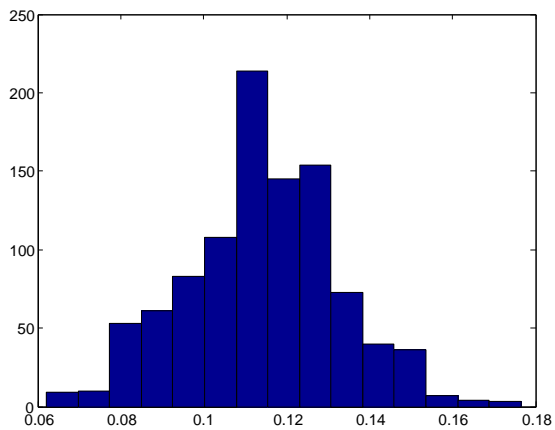
ČSSD, $B = 5000$,
odhad relativními četnostmi



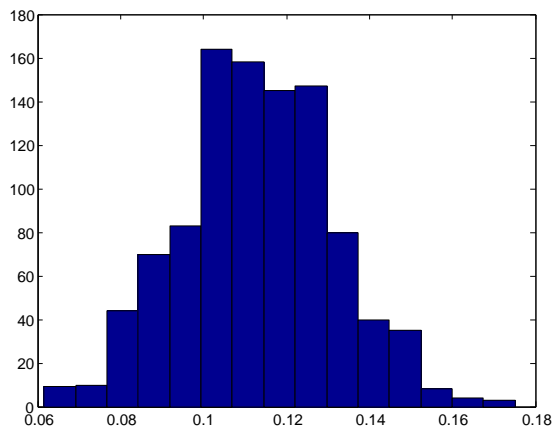
ODS, $B = 1000$,
odhad relativními četnostmi



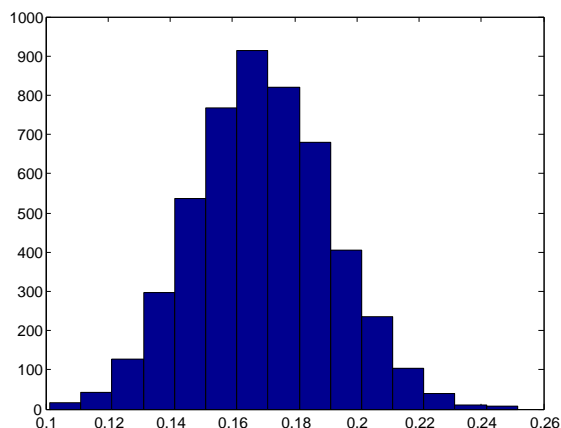
ODS, $B = 5000$,
pesimistický gradientní odhad



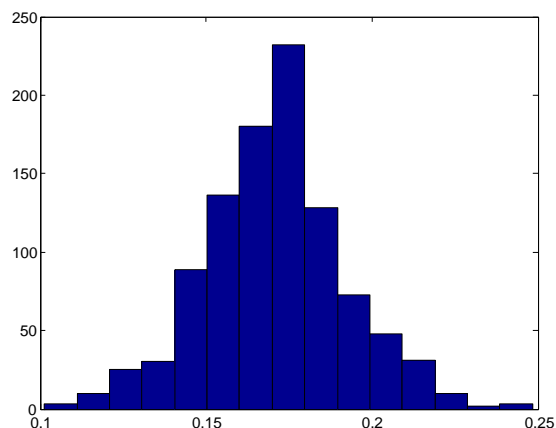
KSČM, $B = 1000$,
odhad relativními četnostmi



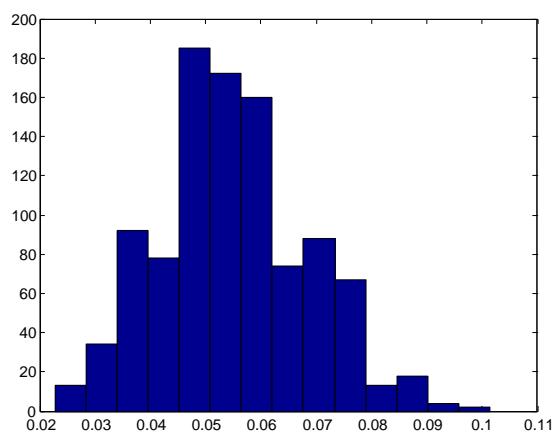
KSČM, $B = 1000$,
pesimistický gradientní odhad



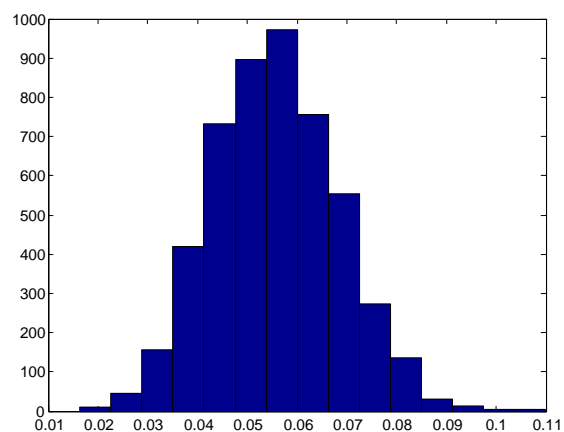
KDU-ČSL, $B = 5000$,
odhad relativními četnostmi



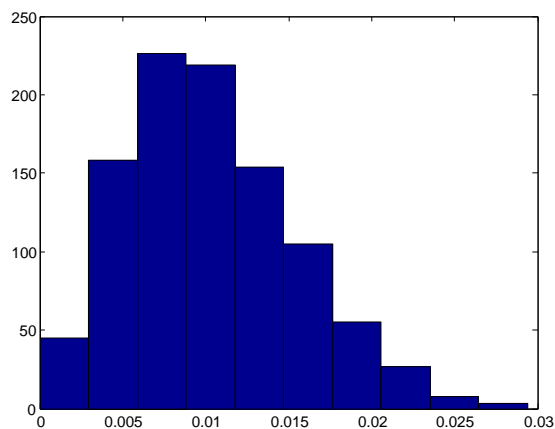
KDU-ČSL, $B = 1000$,
odhad relativními četnostmi



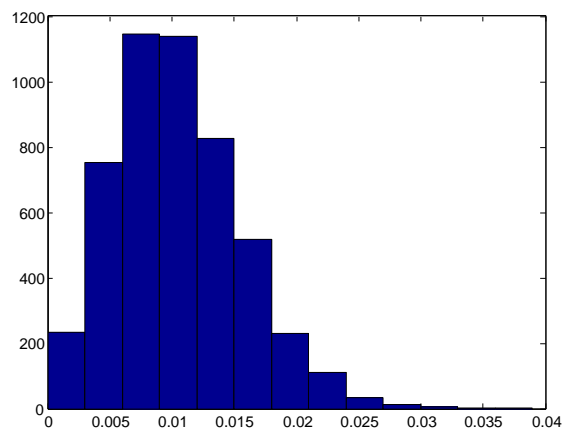
SZ, $B = 1000$,
odhad relativními četnostmi



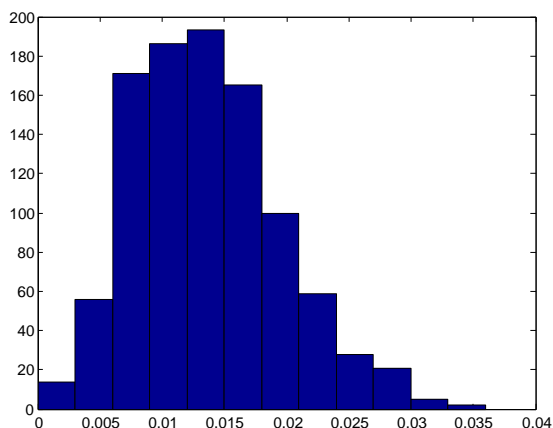
SZ, $B = 5000$,
pesimistický gradientní odhad



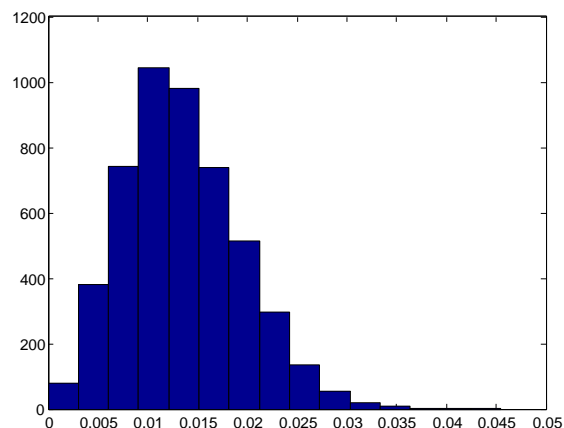
NEZÁVISLÍ, $B = 1000$,
odhad relativními četnostmi



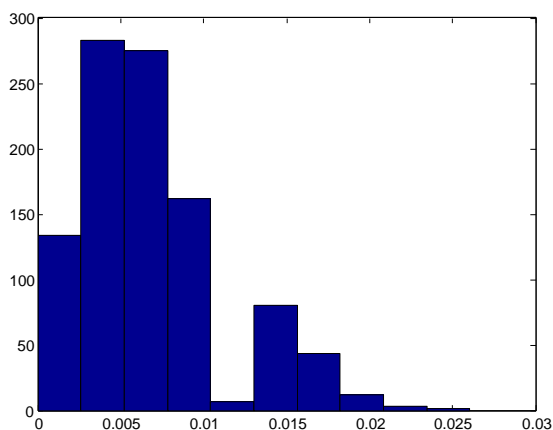
NEZÁVISLÍ, $B = 5000$,
pesimistický gradientní odhad



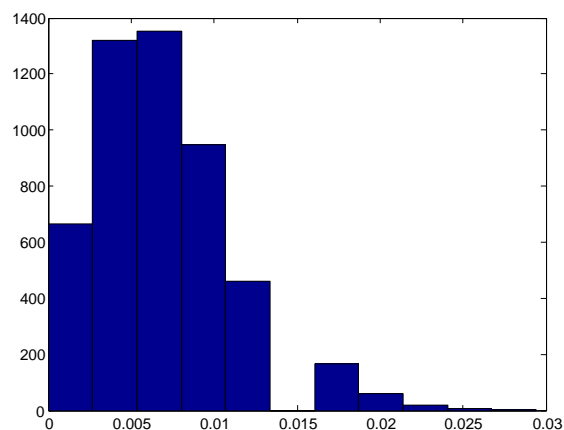
Moravané, $B = 1000$,
odhad relativními četnostmi



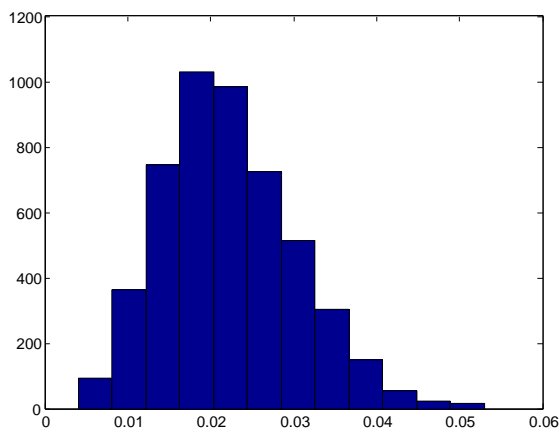
Moravané, $B = 5000$,
pesimistický gradientní odhad



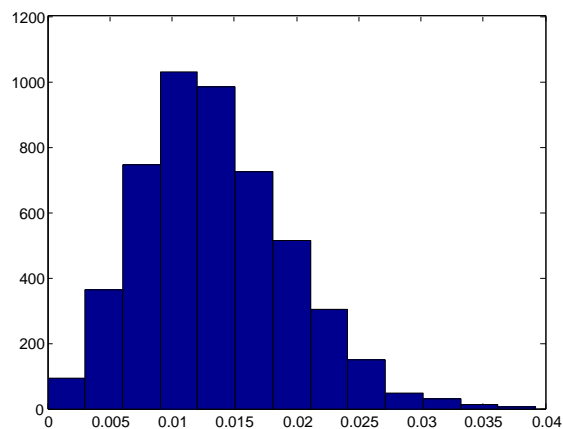
SNK-ED, $B = 1000$,
pesimistický gradientní odhad



SNK-ED, $B = 5000$,
odhad relativními četnostmi



ostatní, $B = 5000$,
pesimistický gradientní odhad



ostatní, $B = 5000$,
odhad relativními četnostmi

všech bootstrapových intervalových odhadů jsou vyšší než horní meze převzatých intervalů, a to o 0,0061 až 0,0183. Jedná se o ty elementární náhodné jevy, u nichž byly pozorované absolutní četnosti nižší než 5, a tedy nebyly splněny předpoklady asymptotické metody použité pro výpočet převzatých intervalových odhadů. Lze se tedy domnívat, že metoda bootstrap v těchto případech přináší „spolehlivější“ intervalové odhady.

Srovnáme-li převzaté i vypočítané intervalové odhady se skutečným výsledkem voleb, zjistíme, že pouze 5 z 9 převzatých intervalových odhadů obsahuje skutečnou hodnotu volebního výsledku, u elementárního náhodného jevu ODS byl volební výsledek nižší, naopak u jevů KDU-ČSL, SNK-ED a ostatní byl vyšší. U jevů ODS, KDU-ČSL a ostatní není skutečný volební výsledek ani prvkem žádného z příslušných bootstrapových intervalových odhadů, zatímco u jevu SNK-ED je prvkem všech čtyř příslušných bootstrapových intervalových odhadů. Dokázali jsme tedy odhadnout volební výsledek úspěšněji než metoda použitá v praxi, i tak ale pouze u šesti elementárních náhodných jevů z devíti. Příčinu musíme hledat zejména v nízkém rozsahu pozorovaného náhodného výběru, 306 respondentů volebního modelu se ukazuje jako zcela nedostatečné množství.

Závěrem lze říci, že metoda bootstrap se ukázala pro odhad pravděpodobnostní funkce volebního modelu jako srovnatelná nebo mírně lepší než v praxi používaná metoda a lze ji doporučit. Zlepšení intervalových odhadů lze dosáhnout zvýšením rozsahu pozorovaného náhodného výběru, zvýšením počtu bootstrapových výběrů a použitím pesimistického gradientního odhadu.

Kapitola 11

Závěr

Připomeneme největší výhody metody bootstrap, jak je uvádí [2]. Umožňuje odhadnout přesnost odhadu, získaného třeba i velmi složitým postupem, užitím výkonu a rychlosti počítače. Její použití zbavuje nutnosti studovat do hloubky teorii a složitě odvozovat přesné vztahy, navíc přináší řešení i v případech, kdy analytické odvození neznáme. Může být použita parametricky i neparametricky. Odhady odvozené neparametrickým bootstrapem jsou pro dostatečně rozsáhlé výběry přesné bez ohledu na pozorované rozdělení pravděpodobnosti. Přitom často jsou dostatečně přesné již pro velmi malé rozsahy. Nicméně základní bootstrap je jen hrubým odhadem přesnosti, který by měl být použit pouze tehdy, když není možné realizovat rozsáhlejší výpočty. Přednost dáváme vždy zkonstruování některého z konfidenčních intervalů pro daný odhad.

Jak se uvádí v [7], metoda pesimistického gradientního odhadu pravděpodobnostní funkce diskrétního rozdělení pravděpodobnosti kategoriální veličiny již na řadě úloh ukázala svoji použitelnost v kategoriální analýze. Až dosud byla ale používána pouze pro bodové odhady. My jsme nyní na příkladech demonstrovali, že ve spojení s metodou bootstrap umožňuje také získat uspokojivé intervalové odhady pravděpodobnostní funkce.

Naším úkolem bylo také otestovat možnosti využití softwaru Shine bootstrap. Jeho použití k výpočtu pesimistického gradientního odhadu pomocí kvadratické kvazinormy lze rozhodně doporučit, bohužel však neumožňuje výpočet kvantilů, a tak je při konstruování intervalových odhadů třeba používat i další software.

Shine bootstrap také umožňuje výpočet maximálně věrohodných bodových odhadů všech tří parametrů tříparametrického Weibullova rozdělení pravděpodobnosti. Je známo, že je-li jeho prahový parametr nenulový, je jeho hodnotu velmi obtížné odhadnout. Shine bootstrap má implementován řešič NOMAD, který funguje na principu ortogonálního prohledávání bez využití derivací. Bohužel se však ukázalo, že spojením takto získaných bodových odhadů parametrů Weibullova rozdělení s metodou bootstrap získáme příliš široké a v případě prahového parametru také příliš vychýlené intervalové odhady, proto tento postup nelze doporučit. Pro získání uspokojivých intervalových odhadů bude zřejmě třeba vyzkoušet implementaci jiných metod výpočtu bodových odhadů parametrů Weibullova rozdělení.

Použité zdroje

- [1] DAVISON, A.C., HINKLEY, D.V. *Bootstrap Methods and Their Applications*. Cambridge: Cambridge University Press, 2003. ISBN 0-521-57471-4.
- [2] EFRON, B., TIBSHIRANI, R.J. *An Introduction to the Bootstrap*. [New York]: Chapman & Hall, 1993. 436 s. ISBN 0-412-04231-2.
- [3] HIGGINS, J.J. *An Introduction to Modern Nonparametric Statistics*. [USA]: Brooks/Cole, 2004. 366 s. ISBN 0-534-38775-6.
- [4] CHERNICK, M.R. *Bootstrap methods: a guide for practitioners and researchers*. New Jersey: John Wiley & Sons, 2008. ISBN 0-471-75621-0.
- [5] KARPÍŠEK, Z. *Matematika IV: Statistika a pravděpodobnost*. 2. doplněné vydání. Brno: Akademické nakladatelství CERM, 2003. 170 s. Skriptum. ISBN 80-214-2522-9.
- [6] KARPÍŠEK, Z., NERADOVÁ, V. *Estimation of Categorical Variable Probability Distribution (Odhad rozdělení pravděpodobnosti kategoriální veličiny)*. Bratislava: 7th International Conference APLIMAT 2008, 5.–8. 2. 2008. Book of abstracts, s. 101. ISBN 978-80-89313-02-0. Sborník, s. 1145–1154. ISBN 978-80-89313-03-7.
- [7] KARPÍŠEK, Z., NERADOVÁ, V., ŽAMPACHOVÁ, E. *A Contribution to the Estimation of Discrete Probability Distribution*. Brno: 14th International Conference on Soft Computing MENDEL 2008, 18.–20. 6. 2008. Sborník, s. 287–292. ISBN 978-80-214-3675-6.
- [8] NERADOVÁ, V. *Progresivní metody odhadů neznámých rozdělení pravděpodobnosti*. Brno: Vysoké učení technické v Brně, Fakulta strojního inženýrství, 2007. Vedoucí diplomové práce doc. RNDr. Zdeněk Karpíšek, CSc.
- [9] *Volby do zastupitelstev krajů konané dne 17. – 18. 10. 2008. Výsledky hlasování za krajské zastupitelstvo. Jihomoravský kraj* [online]. Praha: Český statistický úřad, 2008. URL: <<http://www.volby.cz/pls/kz2008/kz21?xjazyk=CZ&xdatum=20081017&xkraj=10>>
- [10] *Výzkum veřejného mínění. Krajské volby 2008. Jihomoravský kraj* [online]. Praha: SC&C spol. s r. o., 2008-10-8. 47 s. URL: <<http://img1.ct24.cz/multimedia/documents/5/455/45447.doc>>
- [11] *Wikipedia* [online]. URL: <<http://www.wikipedia.com>>