

NMSA331 Matematická statistika I

POZNÁMKY K PŘEDNÁŠCE

Michal Kulich, Marek Omelka

Naposledy upraveno dne 7. ledna 2025.



matfyz

Katedra pravděpodobnosti a matematické statistiky
Matematicko-fyzikální fakulta University Karlovy

Tento učební text představuje modifikaci učebního textu, který připravil **doc. Michal Kulich, Ph.D.** pro kurs *NMSA331 Matematická statistika 1* v rámci bakalářského studia programu *Obecná matematika* na MFF UK. Text vznikl zejména v letech 2015–2022 a následně se již pouze opravovaly objevené překlepy. Nejedná se o plnohodnotnou učebnici ani skriptu, protože zde chybí některé příklady a není zde obsažena látka probíraná na cvičení. Na druhou stranu některé poznatky a poznámky uvedené v tomto textu nebyly probrány na přednášce. Při přípravě na zkoušku je nutné tento text doplnit poznámkami z přednášek a cvičení.

Odkazy na potřebné definice, věty a tvrzení z teorie pravděpodobnosti (začínající písmenem P) se týkají příručky *Základy teorie pravděpodobnosti pro předmět Matematická statistika 1*, která je k dispozici na webových stránkách předmětu NMSA331. Např. tvrzení P.2.2 nebo definici P.6.1 lze najít ve 2., resp. 6. kapitole zmíněné příručky.

Velké poděkování patří také prof. RNDr. Jiřímu Andělovi, DrSc. a doc. RNDr. Karlu Zvárovi, CSc. za pečlivé pročtení poznámek a pomoc s odstraněním řady drobných chyb a nepřesností v prvních verzích tohoto textu. V pozdějších verzích textů pak mnoho nepřesností vychytil doc. Stanislav Nagy, Ph.D.. Dále pak děkujeme všem studentům a studentkám, kteří upozornili na různé nedostatky a překlepy.

OBSAH

ZNAČENÍ	7
1. VYBRANÉ ASYMPTOTICKÉ VÝSLEDKY	10
1.1. Konvergence náhodných vektorů	10
1.2. Zákon velkých čísel	14
1.3. Centrální limitní věta	14
2. NÁHODNÝ VÝBĚR	17
2.1. Definice náhodného výběru	17
2.2. Statistiky	18
2.2.1. Vlastnosti výběrového průměru	18
2.2.2. Relativní četnost	19
2.2.3. Vlastnosti výběrového rozptylu	20
2.3. Uspořádaný náhodný výběr	28
2.4. Transformace ve statistice	35
2.4.1. Transformace pozorování a její vliv na parametry	35
2.4.2. Transformace stabilizující (asymptotický) rozptyl	36
2.4.3. Standardizace	37
3. ODHADOVÁNÍ PARAMETRŮ	40
3.1. Bodový odhad	40
3.2. Volba parametru	43
3.2.1. Kvantitativní data	43
3.2.2. Kategoriální data	44
3.2.3. Binární data	44
3.2.4. Volba parametru v závislosti na typu dat	44
3.3. Momentová metoda	45
3.4. Intervalový odhad	49
3.4.1. Definice	49
3.4.2. Konstrukce intervalových odhadů	52
3.5. Empirické odhady	57
3.5.1. Empirická distribuční funkce	57
3.5.2. Idea empirických odhadů	58
3.5.3. Empirické odhady momentů	58
3.5.4. Empirické odhady kvantilů	60
3.5.5. Empirické odhady pro náhodné vektory	65

4. PRINCIPY TESTOVÁNÍ HYPOTÉZ	69
4.1. Základní pojmy a definice	69
4.2. Hladina a síla testu	71
4.2.1. Hladina testu	72
4.2.2. Síla testu	74
4.2.3. Volba kritického oboru	78
4.3. P-hodnota	85
4.3.1. Výpočet p-hodnoty pro jednostranný kritický obor	86
4.3.2. Výpočet p-hodnoty pro oboustranný kritický obor	88
4.3.3. Rozdělení p-hodnoty za nulové hypotézy	90
4.4. Dualita intervalových odhadů a testování hypotéz	91
5. JEDNOVÝBĚROVÉ A PÁROVÉ PROBLÉMY PRO KVANTITATIVNÍ DATA	94
5.1. Jednovýběrový Kolmogorovův-Smirnovův test	94
5.2. Jednovýběrový t -test	99
5.3. Jednovýběrový znaménkový test	100
5.4. Jednovýběrový Wilcoxonův test	102
5.5. Jednovýběrový χ^2 -test na rozptyl	106
5.6. Párové testy	107
5.7. Párový t -test	108
5.8. Párový znaménkový test	109
5.9. Párový Wilcoxonův test	111
6. DVOUVÝBĚROVÉ PROBLÉMY PRO KVANTITATIVNÍ DATA	114
6.1. Dvouvýběrový Kolmogorovův-Smirnovův test	115
6.2. Dvouvýběrový t -test bez předpokladu shody rozptylů	116
6.3. Dvouvýběrový t -test za předpokladu shody rozptylů	119
6.4. Dvouvýběrový Wilcoxonův test	123
6.5. Dvouvýběrový F -test shody rozptylů	128
7. JEDNOVÝBĚROVÉ A DVOUVÝBĚROVÉ PROBLÉMY PRO BINÁRNÍ DATA	131
7.1. Jednovýběrový problém	131
7.1.1. Clopperova-Pearsonova metoda	131
7.1.2. Klasická asymptotická metoda	133
7.1.3. Wilsonova metoda	133
7.1.4. Logitová metoda	134
7.2. Dvouvýběrový problém	135
7.2.1. Rozdíly pravděpodobností, nárůst rizika	136
7.2.2. Podíly pravděpodobností, relativní riziko	138
7.2.3. Poměr šancí	139
8. MULTINOMICKÉ ROZDĚLENÍ A KONTINGENČNÍ TABULKY	142
8.1. Multinomické rozdělení	142
8.1.1. Odhady parametrů multinomického rozdělení	145

8.1.2.	χ^2 -test dobré shody pro multinomické rozdělení	146
8.1.3.	χ^2 -test dobré shody pro multinomické rozdělení s odhadnutými parametry	148
8.2.	Kontingenční tabulky	150
8.2.1.	Kontingenční tabulky 2×2	154
8.2.2.	Kontingenční tabulky $2 \times K$	155
9.	K-VÝBĚROVÝ PROBLÉM PRO KVANTITATIVNÍ DATA	159
9.1.	Analýza rozptylu (jednoduché třídění)	159
9.2.	Mnohonásobná porovnávání	167
9.2.1.	Bonferroniho metoda	168
9.2.2.	Tukeyova metoda	169
9.3.	Kruskalův-Wallisův test	170
10.	KORELAČNÍ ANALÝZA	175
10.1.	Pearsonův korelační koeficient	175
10.1.1.	Testování hypotézy nezávislosti	176
10.1.2.	Testování obecné hodnoty korelačního koeficientu a interval spolehlivosti pro ρ	177
10.2.	Spearmanův korelační koeficient	178
A.	APPENDIX	182
	APPENDIX	182
A.1.	Rozdělení χ^2 a t	182
A.2.	Idempotentní matice	182
A.3.	Rozdělení kvadratických forem	183
A.4.	Transformace náhodné veličiny její distribuční funkcí	184
A.5.	Gama funkce a beta funkce	184

ZNAČENÍ

\mathbf{a}^\top	transpozice vektoru \mathbf{a}
$\mathbf{a}^{\otimes 2}$	$\mathbf{a}\mathbf{a}^\top$
$\ \mathbf{a}\ $	eukleidovská norma vektoru \mathbf{a}
\xrightarrow{P}	konvergence v pravděpodobnosti
$\xrightarrow{\text{s.j.}}$	konvergence skoro jistě
\xrightarrow{d}	konvergence v distribuci
$X \sim \mathcal{L}$	X má přesné rozdělení \mathcal{L}
$X \overset{\text{as.}}{\sim} \mathcal{L}$	X má přibližně (asymptoticky) rozdělení \mathcal{L}
α	hladina testu
$\beta_n(F), \beta_n(\theta)$	síla testu, silofunkce
γ_3	šikmost náhodné veličiny
$\widehat{\gamma}_3$	empirická šikmost
γ_4	špičatost náhodné veličiny
$\widehat{\gamma}_4$	empirická špičatost
Θ	parametrický prostor
Θ_0	nulová hypotéza
Θ_1	alternativa
λ	Lebesgueova míra na \mathbb{R}
μ_S	čítací míra na nejvýše spočetné množině S
μ_k	k -tý centrální moment náhodné veličiny
$\widehat{\mu}_k$	empirický odhad k -tého centrálního momentu
μ'_k	k -tý moment náhodné veličiny
$\widehat{\mu}'_k$	empirický odhad k -tého momentu
$\rho(X, Y)$	korelační koeficient náhodných veličin X a Y
$\widehat{\rho}_{jm}$	výběrový korelační koeficient j -té a m -té složky náh. vektoru
σ_X	směrodatná odchylka náhodné veličiny X
σ_X^2	rozptyl náhodné veličiny X
$\widehat{\sigma}_n^2$	empirický odhad rozptylu
$\widehat{\Sigma}_n$	výběrová rozptylová matice
φ	hustota normovaného normálního rozdělení
Φ	distribuční funkce normovaného normálního rozdělení

$\chi_f^2(\alpha)$	α -kvantil rozdělení χ_f^2
Ω	prostor elementárních jevů
$\mathbb{1}_B$	indikátor množiny B
$\mathbf{1}_n$	sloupcový vektor jedniček délky n
\mathcal{A}	σ -algebra náhodných jevů na Ω
\mathcal{B}_0	borelovská σ -algebra na \mathbb{R}
\mathcal{B}_0^n	borelovská σ -algebra na \mathbb{R}^n
$C, C(\alpha)$	kritický obor testu
$c_L(\alpha), c_U(\alpha)$	kritické hodnoty
$\text{cor}(X, Y)$	korelační koeficient náhodných veličin X a Y
$\text{cor}(\mathbf{X}, \mathbf{Y})$	korelační matice náhodných vektorů \mathbf{X} a \mathbf{Y}
$\text{cov}(X_1, X_2)$	kovariance náhodných veličin X_1 a X_2
$\text{cov}(\mathbf{X}_1, \mathbf{X}_2)$	kovarianční matice náhodných vektorů \mathbf{X}_1 a \mathbf{X}_2
$\text{diag}(\mathbf{a})$	diagonální matice obsahující složky vektoru \mathbf{a} na diagonále
$E X$	střední hodnota náhodné veličiny (vektoru) X
\mathcal{F}	pravděpodobnostní model pro pozorovaná data
\mathcal{F}_0	rozdělení splňující nulovou hypotézu
\mathcal{F}_1	rozdělení splňující alternativu
f_X	hustota náhodné veličiny (vektoru) X
F_X	distribuční funkce náhodné veličiny (vektoru) X
F_X^{-1}	kvantilová funkce náhodné veličiny X
\widehat{F}_n	empirická distribuční funkce
$F_{m,n}(\alpha)$	α -kvantil rozdělení $F_{m,n}$
H_0	nulová hypotéza
H_1	alternativa
$\mathbb{1}_n$	jednotková matice $n \times n$
\mathcal{L}^p	množina náhodných veličin na (Ω, \mathcal{A}, P) s konečným p -tým absolutním momentem
\mathcal{L}_+^2	množina náhodných veličin na (Ω, \mathcal{A}, P) s konečným a nenulovým rozptylem
$\mathcal{L}(X)$	rozdělení náhodné veličiny (vektoru) X
m_X	medián náhodné veličiny X
\widehat{m}_n	výběrový medián
MSE	střední čtvercová odchylka odhadu

P	pravděpodobnost
P_X	rozdělení náhodné veličiny X , její indukovaná míra na výběrovém prostoru
P_θ	rozdělení dat při hodnotě parametru θ
$h(\mathbb{A})$	hodnota matice \mathbb{A}
\mathbb{R}	množina reálných čísel
R_i	pořadí i -tého pozorování
SE	směrodatná chyba odhadu
S_n^2	výběrový rozptyl
S_{jm}	výběrová kovariance j -té a m -té složky náh. vektoru
S_X	nosič rozdělení náhodné veličiny X
$t_f(\alpha)$	α -kvantil rozdělení t_f
$\text{tr}(\mathbb{A})$	stopa matice \mathbb{A}
$u_X(\alpha)$	α -kvantil náhodné veličiny X
u_α	α -kvantil rozdělení $N(0, 1)$
$\hat{u}_n(\alpha)$	výběrový α -kvantil
$\text{var } X$	rozptyl náhodné veličiny X
$\text{var } \mathbf{X}$	rozptylová (varianční) matice náhodného vektoru \mathbf{X}
\mathcal{X}	výběrový prostor
$X_{(k)}$	k -tá pořádková statistika
\bar{X}_n	výběrový průměr náhodného výběru X_1, \dots, X_n

1. VYBRANÉ ASYMPTOTICKÉ VÝSLEDKY

Mějme posloupnost k -rozměrných náhodných vektorů X_1, X_2, X_3, \dots , kde vektor $X_i = (X_{i1}, \dots, X_{ik})^T$ je definován na $(\Omega_i, \mathcal{A}_i, P_i)$.

1.1. KONVERGENCE NÁHODNÝCH VEKTORŮ

Nechť $\|a\|$ značí eukleidovskou normu vektoru a , tj. $\|a\| = \sqrt{a^T a}$.

Definice 1.1 (konvergence skoro jistě) Říkáme, že posloupnost náhodných vektorů $\{X_n\}_{n=1}^\infty$ *konverguje skoro jistě* k náhodnému vektoru X pro $n \rightarrow \infty$ právě když

$$P(\omega : \lim_{n \rightarrow \infty} \|X_n(\omega) - X(\omega)\| = 0) = 1.$$

Konvergenci skoro jistě značíme $X_n \xrightarrow[n \rightarrow \infty]{\text{s.j.}} X$.

Definice 1.2 (konvergence v pravděpodobnosti) Říkáme, že posloupnost náhodných vektorů $\{X_n\}_{n=1}^\infty$ *konverguje v pravděpodobnosti* k náhodnému vektoru X pro $n \rightarrow \infty$ právě když

$$\forall \varepsilon > 0 : \lim_{n \rightarrow \infty} P(\omega : \|X_n(\omega) - X(\omega)\| > \varepsilon) = 0.$$

Konvergenci v pravděpodobnosti značíme $X_n \xrightarrow[n \rightarrow \infty]{P} X$.

Poznámka.

- Pro $k = 1$ odpovídají definice 1.1 a 1.2 příslušným definicím z předmětu NMSA202 (*Pravděpodobnost a matematická statistika*).
- S využitím nerovností

$$\max_{j \in \{1, \dots, k\}} |a_j| \leq \|a\| \leq \sqrt{k} \max_{j \in \{1, \dots, k\}} |a_j|$$

Ize alternativně definovat konvergenci skoro jistě a v pravděpodobnosti pro náhodné vektory po složkách, tj.

$$X_n \xrightarrow[n \rightarrow \infty]{\text{s.j.}} X \Leftrightarrow X_{nj} \xrightarrow[n \rightarrow \infty]{\text{s.j.}} X_j, \forall j \in \{1, \dots, k\},$$

$$X_n \xrightarrow[n \rightarrow \infty]{P} X \Leftrightarrow X_{nj} \xrightarrow[n \rightarrow \infty]{P} X_j, \forall j \in \{1, \dots, k\}.$$

- Aby definice 1.1 dávala smysl, musí být všechny náhodné vektory definované na stejném pravděpodobnostním prostoru (Ω, \mathcal{A}, P) . U definice 1.2 to není nutné, pokud limitní náhodný vektor X je rovný konstantě skoro jistě.

V následujícím budeme značit F_{X_n} distribuční funkci náhodného vektoru X_n , tj.

$$F_{X_n}(x) = P(X_n \leq x).$$

Podobně F_X bude distribuční funkce náhodného vektoru X .

Definice 1.3 (konvergence v distribuci) Říkáme, že posloupnost $\{X_n\}_{n=1}^{\infty}$ konverguje v distribuci k náhodnému vektoru X pro $n \rightarrow \infty$ právě když

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$$

v každém bodě x , v němž je $F_X(x)$ spojitá. Konvergenci v distribuci značíme $X_n \xrightarrow[n \rightarrow \infty]{d} X$.

Příklad. Necht' X_1, \dots, X_n jsou nezávislé stejně rozdělené náhodné veličiny s rozptylem $\text{var } X_1 \in (0, \infty)$. Potom z předmětu NMSA 202 (Pravděpodobnost a matematická statistika) víme, že $\frac{\sqrt{n}(\bar{X}_n - EX_1)}{\sqrt{\text{var } X_1}} \xrightarrow[n \rightarrow \infty]{d} Z$, kde $Z \sim N(0, 1)$. Jelikož distribuční funkce Φ náhodné veličiny s rozdělením $N(0, 1)$ je spojitá \mathbb{R} , tak dle definice konvergence v distribuci máme

$$\lim_{n \rightarrow \infty} P\left(\frac{\sqrt{n}(\bar{X}_n - EX_1)}{\sqrt{\text{var } X_1}} \leq x\right) = \Phi(x), \quad \forall x \in \mathbb{R}.$$

Příklad. Necht' náhodná veličina U má rovnoměrné rozdělení na intervalu $(0, 1)$ a položme $X_n = U^n$. S využitím definice 1.3 dokažte, že $X_n \xrightarrow[n \rightarrow \infty]{d} 0$.

Poznámka.

- Pro konvergence v distribuci je podstatné pouze to, co se děje s rozdělením náhodného vektoru. Označíme-li tedy $\mathcal{L}(X_n)$ rozdělení náhodného vektoru X_n (z angl. Law), pak konvergenci v distribuci můžeme zapisovat také jako $\mathcal{L}(X_n) \rightarrow \mathcal{L}(X)$ pro $n \rightarrow \infty$. Tento zápis pak čteme, že „rozdělení X_n konverguje k rozdělení X “. Můžeme také říkat, že X_n má asymptotické (limitní) rozdělení F_X a psát $X_n \overset{\text{as.}}{\sim} \mathcal{L}(X)$.
- Pro konvergenci v distribuci nepotřebujeme, aby náhodné vektory byly definovány na stejném pravděpodobnostním prostoru.
- Na rozdíl od konvergence skoro jistě a v pravděpodobnosti, tak konvergenci v distribuci nelze ekvivalentně definovat po složkách.

Tvrzení 1.1

$$(i) \quad X_n \xrightarrow[n \rightarrow \infty]{\text{s.j.}} X \Rightarrow X_n \xrightarrow[n \rightarrow \infty]{P} X$$

$$(ii) \mathbf{X}_n \xrightarrow[n \rightarrow \infty]{P} \mathbf{X} \Rightarrow \mathbf{X}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{X}$$

Poznámka. Opačné implikace neplatí. Nicméně pokud náhodné vektory konvergují v distribuci ke konstantě, tj. $\mathbf{X}_n \xrightarrow{d} c$ (kde $c \in \mathbb{R}^k$), pak platí $\mathbf{X}_n \xrightarrow{P} c$.

Následující věta říká, že spojitá transformace zachovává všechny výše uvedené druhy konvergencí.

Tvrzení 1.2 (Věta o spojitě transformaci) Necht' $\mathbf{X}, \mathbf{X}_1, \mathbf{X}_2, \dots$ jsou náhodné vektory a funkce $g: \mathbb{R}^k \rightarrow \mathbb{R}^m$ je spojitá na množině C takové, že $P(\mathbf{X} \in C) = 1$. Potom:

$$(i) \mathbf{X}_n \xrightarrow[n \rightarrow \infty]{s.j.} \mathbf{X} \Rightarrow g(\mathbf{X}_n) \xrightarrow[n \rightarrow \infty]{s.j.} g(\mathbf{X});$$

$$(ii) \mathbf{X}_n \xrightarrow[n \rightarrow \infty]{P} \mathbf{X} \Rightarrow g(\mathbf{X}_n) \xrightarrow[n \rightarrow \infty]{P} g(\mathbf{X});$$

$$(iii) \mathbf{X}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{X} \Rightarrow g(\mathbf{X}_n) \xrightarrow[n \rightarrow \infty]{d} g(\mathbf{X}).$$

Důkaz. Část (i)

$$\begin{aligned} & P \left[\omega : \lim_{n \rightarrow \infty} \|g(\mathbf{X}_n(\omega)) - g(\mathbf{X}(\omega))\| = 0 \right] \\ & \geq P \left[\omega : \lim_{n \rightarrow \infty} \|g(\mathbf{X}_n(\omega)) - g(\mathbf{X}(\omega))\| = 0, \mathbf{X}(\omega) \in C \right] \\ & = P \left[\omega : \lim_{n \rightarrow \infty} \|\mathbf{X}_n(\omega) - \mathbf{X}(\omega)\| = 0, \mathbf{X}(\omega) \in C \right] = 1, \end{aligned}$$

kde jsme využili toho, že g je spojitá na množině C a $P(\mathbf{X} \in C) = 1$.

Část (ii) Necht' $\varepsilon > 0$. Potom pro všechna $\delta > 0$

$$\begin{aligned} & P \left[\omega : \|g(\mathbf{X}_n(\omega)) - g(\mathbf{X}(\omega))\| > \varepsilon \right] \\ & \leq P \left[\|g(\mathbf{X}_n) - g(\mathbf{X})\| > \varepsilon, \|\mathbf{X}_n - \mathbf{X}\| \leq \delta \right] + P \left[\|\mathbf{X}_n - \mathbf{X}\| > \delta \right] \\ & \leq P \left[\mathbf{X} \in B^\delta \right] + \underbrace{P \left[\|\mathbf{X}_n - \mathbf{X}\| > \delta \right]}_{\rightarrow 0, \forall \delta > 0}, \end{aligned}$$

kde $B^\delta = \{x \in \mathbb{R}^k; \exists y \in \mathbb{R}^k : \|x - y\| \leq \delta, \|g(x) - g(y)\| > \varepsilon\}$. Dále

$$\begin{aligned} P[\mathbf{X} \in B^\delta] & \leq P[\mathbf{X} \in B^\delta, \mathbf{X} \in C] + P[\mathbf{X} \in B^\delta, \mathbf{X} \notin C] \\ & = P[\mathbf{X} \in B^\delta \cap C] + 0. \end{aligned}$$

Dále ze spojitosti funkce g na množině C platí, že $B^\delta \cap C \rightarrow \emptyset$ pro $\delta \searrow 0$. Tudíž pravděpodobnost $P[\mathbf{X} \in B^\delta \cap C]$ můžeme udělat (pro všechna dostatečně velká n) libovolně malou tím, že volíme dostatečně malé δ .

Část (iii) Viz důkaz věty 13.6 in [Lachout \(2004\)](#).

□

Ve statistice budeme často používat následující větu, která je zobecněním věty 4.14 z [Dupač and Hušková \(1999\)](#).

Tvrzení 1.3 (Cramérová-Sluckého věta) Nechť $\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{X}$, $\mathbb{A}_n \xrightarrow[n \rightarrow \infty]{P} \mathbb{A}$ a $\mathbf{B}_n \xrightarrow[n \rightarrow \infty]{P} \mathbf{b}$, kde \mathbf{X}_n a \mathbf{X} jsou k -rozměrné náhodné vektory, \mathbb{A}_n je náhodná matice o dimenzích $m \times k$, \mathbb{A} je matice konstant o dimenzích $m \times k$, \mathbf{B}_n jsou m -rozměrné náhodné vektory a \mathbf{b} je m -rozměrný vektor konstant, pak

$$\mathbb{A}_n \mathbf{X}_n + \mathbf{B}_n \xrightarrow[n \rightarrow \infty]{d} \mathbb{A} \mathbf{X} + \mathbf{b}.$$

Poznámka. Cramérově-Sluckého větě se často říká pouze Sluckého věta.

Tvrzení 1.4 Nechť $a_n(\mathbf{X}_n - \boldsymbol{\mu}) \xrightarrow[n \rightarrow \infty]{d} \mathbf{X}$, kde $a_n > 0$ je posloupnost reálných čísel splňující $a_n \rightarrow \infty$ pro $n \rightarrow \infty$ a $\boldsymbol{\mu}$ je vektor konstant. Pak $\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\mu}$.

Důkaz. Vzhledem k poznámce za definicí konvergence v pravděpodobnosti (definice 1.2) stačí dokázat, že pro všechna $j \in \{1, \dots, k\}$ platí $X_{nj} \xrightarrow[n \rightarrow \infty]{P} \mu_j$, kde μ_j je j -tá složka vektoru $\boldsymbol{\mu}$. Navíc bez újmy na obecnosti můžeme předpokládat, že $\mu_j = 0$.

Nechť $\varepsilon > 0$ je dáno. Potom pro libovolně malé $\eta > 0$ můžeme najít konečnou konstantu K takovou, že $P(|X_j| \geq K) < \eta$ a zároveň K a $-K$ jsou body spojitosti distribuční funkce X_j . Potom pro všechna n taková, že $a_n \varepsilon > K$, platí

$$P(|X_{nj}| > \varepsilon) = P(a_n |X_{nj}| > a_n \varepsilon) \leq P(a_n |X_{nj}| > K) \leq 1 - F_{nj}(K) + F_{nj}(-K), \quad (1.1)$$

kde F_{nj} je distribuční funkce náhodné veličiny $a_n(X_{nj} - \mu_j)$.

Dále z předpokladu $a_n(\mathbf{X}_n - \boldsymbol{\mu}) \xrightarrow[n \rightarrow \infty]{d} \mathbf{X}$, plyne speciálně že $a_n(X_{nj} - \mu_j) \xrightarrow[n \rightarrow \infty]{d} X_j$. Tudíž z definice konvergence v distribuci

$$\lim_{n \rightarrow \infty} [1 - F_{nj}(K) + F_{nj}(-K)] = 1 - F_{X_j}(K) + F_{X_j}(-K) \leq P(|X_j| \geq K) < \eta. \quad (1.2)$$

Tedy celkem z (1.1) a (1.2) plyne, že pro všechna dostatečně velká n platí

$$P(|X_{nj}| > \varepsilon) < 2\eta.$$

Jelikož η lze vzít libovolně malé, tak odsud plyne, že $\lim_{n \rightarrow \infty} P(|X_{nj}| > \varepsilon) = 0$, čímž je důkaz dokončen. \square

1.2. ZÁKON VELKÝCH ČÍSEL

Tvrzení 1.5 Necht' X_1, X_2, \dots je posloupnost nezávislých stejně rozdělených náhodných vektorů s konečnou střední hodnotou $E X_i = \mu$. Pak platí

$$\overline{X}_n \xrightarrow{\text{s.j.}} \mu.$$

Důkaz. Důkaz plyne použitím Kolmogorovova silného zákona velkých čísel na jednotlivé složky náhodného vektoru. \square

1.3. CENTRÁLNÍ LIMITNÍ VĚTA

Tvrzení 1.6 (centrální limitní věta pro nezávislé stejně rozdělené náhodné vektory) Necht' X_1, X_2, \dots jsou nezávislé a stejně rozdělené náhodné vektory se střední hodnotou $E X_i = \mu$ a konečnou rozptylovou maticí $\text{var } X_i = \Sigma$. Pak platí

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu) = \sqrt{n} (\overline{X}_n - \mu) \xrightarrow[n \rightarrow \infty]{\text{d}} N_k(\mathbf{0}, \Sigma).$$

Poznámka. Neformální zápis tvrzení centrální limitní věty: $\overline{X}_n \stackrel{\text{as.}}{\approx} N_k(\mu, n^{-1}\Sigma)$.

Tvrzení 1.7 (Δ -metoda) Necht' $\{T_n\}_{n=1}^\infty$ splňuje

$$\sqrt{n} (T_n - \mu) \xrightarrow[n \rightarrow \infty]{\text{d}} N_k(\mathbf{0}, \Sigma) \quad (1.3)$$

pro nějaký vektor konstant $\mu \in \mathbb{R}^k$ a matici Σ . Necht' $g : \mathbb{R}^k \rightarrow \mathbb{R}^p$ je funkce, která je spojitě diferencovatelná v nějakém okolí bodu μ . Označme $\mathbb{D}(x) = \frac{\partial g(x)}{\partial x}$. Pak platí

$$\sqrt{n} (g(T_n) - g(\mu)) \xrightarrow[n \rightarrow \infty]{\text{d}} N_p(\mathbf{0}, \mathbb{D}(\mu) \Sigma \mathbb{D}(\mu)^\top).$$

Důkaz. Pro dané $j \in \{1, \dots, p\}$ uvažujme $g_j : \mathbb{R}^k \rightarrow \mathbb{R}$ (tj. j -tou složku zobrazení g). Z předpokladů věty existuje okolí \mathcal{U} bodu μ takové, že g_j má spojitě parciální derivace na tomto okolí. Dále z předpokladu (1.3) a tvrzení 1.4 plyne, že T_n konverguje k μ v pravděpodobnosti, tj. $P(T_n \in \mathcal{U}) \rightarrow 1$ pro $n \rightarrow \infty$. Tedy bez újmy na obecnosti můžeme předpokládat, že $T_n \in \mathcal{U}$. Nyní dle věty o střední hodnotě existuje μ_n^{j*} , které leží mezi T_n a μ takové, že

$$\sqrt{n} (g_j(T_n) - g_j(\mu)) = \nabla g_j(\mu_n^{j*}) \sqrt{n} (T_n - \mu).$$

Nechť \mathbb{D}_n je matice typu $p \times k$, která má v j -tém řádku $\nabla g_j(\mu_n^{j*})$ a všimněme si, že

$$\sqrt{n} (g(T_n) - g(\mu)) = \mathbb{D}_n \sqrt{n} (T_n - \mu). \quad (1.4)$$

Nyní $T_n \xrightarrow[n \rightarrow \infty]{P} \mu$ implikuje, že $\mu_n^{j*} \xrightarrow[n \rightarrow \infty]{P} \mu$. S využitím spojitosti parciálních derivací zobrazení g na \mathcal{U} a tvrzení 1.2(ii) dostáváme

$$\mathbb{D}_n \xrightarrow[n \rightarrow \infty]{P} \mathbb{D}(\mu).$$

Toto společně s (1.4) a tvrzením 1.3 dokončuje důkaz. \square

Poznámka.

- Tvrzení 1.7 budeme nejčastěji používat v jednorozměrném případě. Tj. necht

$$\sqrt{n} (T_n - \mu) \xrightarrow[n \rightarrow \infty]{d} N(0, \sigma^2)$$

a funkce $g : \mathbb{R} \rightarrow \mathbb{R}$ má spojitou derivaci na nějakém okolí bodu μ . Pak platí

$$\sqrt{n}(g(T_n) - g(\mu)) \xrightarrow[n \rightarrow \infty]{d} N(0, [g'(\mu)]^2 \sigma^2). \quad (1.5)$$

Předchozí výsledek se dá pamatovat následovně. S využitím Taylorova rozvoje prvního řádu funkce $g(T_n)$ v bodě μ můžeme přepsat

$$\sqrt{n} (g(T_n) - g(\mu)) = g'(\mu) \sqrt{n} (T_n - \mu) + R_n, \quad (1.6)$$

kde R_n je zbytkový člen. Z předpokladu (1.3) ihned plyne, že

$$g'(\mu) \sqrt{n} (T_n - \mu) \xrightarrow[n \rightarrow \infty]{d} N(0, [g'(\mu)]^2 \sigma^2).$$

Tudíž pokud můžeme zanedbat zbytkový člen R_n (přesněji $R_n \xrightarrow[n \rightarrow \infty]{P} 0$), tak z (1.6) dostáváme (1.5).

Podobně se dá snadno zapamatovat i tvrzení 1.7.

- Jako T_n budeme nejčastěji brát \bar{X}_n , kde $\mathbf{X}_i = (X_{i1}, \dots, X_{ik})^\top$ jsou vhodně zvolené nezávislé stejně rozdělené náhodné vektory. K ověření předpokladu (1.3) pak můžeme využít centrální limitní větu (tvrzení 1.6).
- Všimněme si, že tvrzení 1.7 společně s tvrzením 1.4 implikují, že

$$g(T_n) \xrightarrow[n \rightarrow \infty]{P} g(\mu). \quad (1.7)$$

Pozorný čtenář by však mohl namítnout, že pokud by nás zajímal pouze výsledek (1.7), pak předpoklady tvrzení 1.7 jsou zbytečně silné. Dle věty o spojitě transformaci (tvrzení 1.2) by nám k důkazu konvergence (1.7) stačilo, že funkce g je spojitá v bodě μ .

Přípravné příklady ke zkoušce. Necht' X_1, X_2, \dots jsou nezávislé a stejně rozdělené náhodné veličiny. Označme $\mu = E X_i$ a $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$.

1. Určete k čemu a v jakém smyslu konverguje \bar{X}_n^2 .
2. Necht' navíc platí, že $\mu = 0$ a $\text{var } X_i = 1$. Určete k čemu a v jakém smyslu konvergují následující náhodné veličny: (a) $\sqrt{n}(\bar{X}_n)^2$, (b) $n(\bar{X}_n)^2$ a (c) $n^{3/2}(\bar{X}_n)^2$.

2. NÁHODNÝ VÝBĚR

2.1. DEFINICE NÁHODNÉHO VÝBĚRU

Nechť je dán pravděpodobnostní prostor (Ω, \mathcal{A}, P) .

Definice 2.1 Posloupnost X_1, X_2, \dots, X_n nezávislých stejně rozdělených náhodných vektorů definovaných na (Ω, \mathcal{A}, P) , z nichž každý má distribuční funkci F_X , nazýváme *náhodný výběr z rozdělení F_X* .^{*} Konstantu n nazýváme *rozsah výběru*.[†]

Prvky náhodného výběru mohou být buď reálné náhodné veličiny nebo náhodné vektory (matice apod.). Můžeme je nazývat „pozorování“ nebo „data“. Pro označení náhodného výběru jako celku budeme občas používat značení X .

Poznámka. Distribuční funkci F_X , z níž pozorování X_1, X_2, \dots, X_n pocházejí, neznáme. Chceme použít pozorování k tomu, abychom se o F_X něco potřebného dozvěděli. O distribuční funkci F_X předpokládáme, že patří do nějaké množiny rozdělení \mathcal{F} , které říkáme *model*.

Definice 2.2 *Modelem* pro náhodný výběr X_1, X_2, \dots, X_n rozumíme předem stanovenou množinu rozdělení \mathcal{F} , do níž patří neznámé rozdělení F_X .

Poznámka. Rozdělení F_X je neznámé. Rádi bychom použili pozorovaná data X , abychom určili jeho jisté charakteristiky, které nazýváme *parametry*. Formálně jde o nějakou konstantu (nebo vektor konstant) $\theta_X \in \mathbb{R}^k$, kterou bychom uměli zjistit, kdybychom F_X znali. Hledaný parametr tedy můžeme obecně zapsat ve tvaru $\theta_X \equiv t(F_X)$, kde t je nějaký funkcional.

Příklady (Typy modelů pro reálné náhodné veličiny).

1. Za model \mathcal{F} můžeme např. vzít množinu všech [diskrétních, spojitých] rozdělení na \mathbb{R} s konečnou střední hodnotou [s konečným rozptylem]. Hledané parametry mohou být např. $E X_i$, $\text{var } X_i$, $P[X \leq x] \equiv F_X(x)$ nebo kvantil $F_X^{-1}(\alpha)$. Takový model nazýváme *neparametrický*[‡], neboť není možné popsat všechna rozdělení v \mathcal{F} pomocí konečně mnoha parametrů. Symbolem Θ označujeme množinu všech přípustných hodnot parametru $\theta \equiv t(F)$ pro všechna $F \in \mathcal{F}$.
2. Za model \mathcal{F} můžeme vzít množinu všech rozdělení s hustotami tvaru $f(x; \theta)$ pro $\theta \in \Theta \subseteq \mathbb{R}^p$, kde $f(\cdot; \cdot)$ je známá funkce a θ je neznámá konstanta (např. všechna exponenciální, normální, geometrická rozdělení). Tyto modely nazýváme *parametrické*[§]. V parametrickém modelu lze jakékoli jiné parametry vždy vyjádřit jako funkce θ .

^{*} Angl. *random sample from distribution F_X* [†] Angl. *sample size* [‡] Angl. *non-parametric model*
[§] Angl. *parametric model*

Příklady (Parametrické modely).

- $\mathcal{F} = \{N(\mu, \sigma_0^2), \mu \in \mathbb{R}, \sigma_0^2 \text{ pevně dáno}\}; \theta = \mu, \Theta = \mathbb{R}$.
- $\mathcal{F} = \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}^+\}; \theta = (\mu, \sigma^2)^T, \Theta = \mathbb{R} \times \mathbb{R}^+$.
- $\mathcal{F} = \{\text{Exp}(\lambda), \lambda \in \mathbb{R}^+\}; \theta = \lambda, \Theta = \mathbb{R}^+$.
- $\mathcal{F} = \{\text{Alt}(p), p \in (0, 1)\}; \theta = p, \Theta = (0, 1)$.

Poznámka. Model \mathcal{F} a parametr θ , který nás zajímá, volíme sami. Model vyjadřuje naši apriorní (na datech nezávislou) představu o rozdělení pozorovaných veličin. Volba parametru závisí na otázce, kterou se snažíme zodpovědět pomocí statistické analýzy. Volba modelu a parametru ovlivňuje výběr metody pro analýzu dat (a její výsledky).

2.2. STATISTIKY

Statistická analýza postupuje tak, že se z náhodného výběru počítají veličiny, které obsahují informaci o požadovaných parametrech, a s nimi se dále pracuje. Těmto veličinám se říká statistiky. Uvažujme náhodný výběr $\mathbf{X} = (X_1, X_2, \dots, X_n)$.

Definice 2.3 Pojmem *statistika** nazýváme libovolnou měřitelnou funkci $S(\mathbf{X})$ pozorování z náhodného výběru \mathbf{X} . Statistika je náhodná veličina (náhodný vektor, je-li vícerozměrná).

Statistika nesmí záviset na hodnotách, které neznáme a nepozorujeme. Smí to být pouze funkce dat a známých konstant. Mezi nejčastěji používané statistiky patří výběrový průměr a výběrový rozptyl. Uvažujme nyní výběr reálných náhodných veličin $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$ a zaveďme dvě nejčastěji používané statistiky.

Definice 2.4

- Veličina $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ se nazývá *výběrový průměr*† náhodného výběru \mathbf{X} .
- Pro $n \geq 2$ veličina $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ se nazývá *výběrový rozptyl*‡ náhodného výběru \mathbf{X} .

Výběrový rozptyl nemá smysl počítat z jediného pozorování ($n = 1$); uvažujme-li výběrový rozptyl, automaticky předpokládáme, že $n \geq 2$.

2.2.1. VLASTNOSTI VÝBĚROVÉHO PRŮMĚRU

Uvažujme obecný model $\mathcal{F} = \mathcal{L}^2$. Pracujme tedy s náhodným výběrem \mathbf{X} , jehož složky X_i jsou nezávislé náhodné veličiny s libovolným rozdělením, které má konečné druhé momenty. Označme $\mu \equiv E X_i$ a $\sigma^2 = \text{var } X_i$.

Lemma 2.1

$$\bar{X}_n = \arg \min_{c \in \mathbb{R}} \sum_{i=1}^n (X_i - c)^2.$$

* Angl. *statistic* † Angl. *sample mean* ‡ Angl. *sample variance*

Důkaz. Označme si funkci $f(c) = \sum_{i=1}^n (X_i - c)^2$. Tvrzení lemmatu plyne z toho, že $f'(\bar{X}_n) = 0$ a $f''(c) > 0$ pro všechna $c \in \mathbb{R}$. \square

Výběrový průměr tedy minimalizuje součet čtverců odchylek jednotlivých pozorování od libovolného reálného čísla.

Snadno spočítáme první dva momenty výběrového průměru a prozkoumáme jeho limitní chování při $n \rightarrow \infty$.

Věta 2.2 (Vlastnosti průměru)

- (i) $E \bar{X}_n = \mu$, $\text{var } \bar{X}_n = \frac{\sigma^2}{n}$;
- (ii) $\bar{X}_n \xrightarrow[n \rightarrow \infty]{P} \mu$;
- (iii) $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow[n \rightarrow \infty]{d} N(0, \sigma^2)$ neboli $\bar{X}_n \stackrel{\text{as.}}{\sim} N(\mu, \frac{\sigma^2}{n})$.

Důkaz. (i) plyne z přímého výpočtu. (ii) ze silného zákona velkých čísel (tvrzení 1.5 pro $k = 1$) a (iii) z centrální limitní věty (tvrzení 1.6 pro $k = 1$). \square

Poznámka. Platí-li předpoklad normálního rozdělení, tj. $\mathcal{F} = \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}^+\}$, lze body (i) a (iii) předchozí věty zesílit na

$$\sqrt{n}(\bar{X}_n - \mu) \sim N(0, \sigma^2) \quad \text{neboli} \quad \bar{X}_n \sim N(\mu, \frac{\sigma^2}{n}).$$

Důkaz. Z předpokladů plyne, že náhodný vektor $Z = (X_1 - \mu, \dots, X_n - \mu)^\top$ má nezávislé složky s rozdělením $N(0, \sigma^2)$ a z definice mnohorozměrného normálního rozdělení (viz definice P.6.1) plyne, že $Z \sim N_n(\mathbf{0}, \sigma^2 \mathbb{1}_n)$. Označme si $c = (\frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}})^\top \in \mathbb{R}^n$. Z vlastnosti mnohorozměrného normálního rozdělení plyne, že

$$c^\top Z = \sqrt{n}(\bar{X}_n - \mu) \sim N(0, \sigma^2).$$

\square

2.2.2. RELATIVNÍ ČETNOST

V aplikacích často nabývá náhodná veličina X_i pouze hodnot 0 a 1. Jednička při tom znamená, že v i -tém pokusu nastal nějaký jev B a nula, že nenastal. Označme $p = P(X_i = 1)$. Pak náhodné veličiny X_1, \dots, X_n představují náhodný výběr z *alternativního rozdělení** $\text{Alt}(p)$.

Výběrový průměr \bar{X}_n je podílem počtu pozorování, při nichž jev B nastal, a celkového počtu pozorování n . Nazýváme jej (*empirická*) *relativní četnost*[†] jevu B . Pro relativní četnost \bar{X}_n pochopitelně platí věta 2.2. Uveďme si ji znovu v podobě specializované na tento případ a přidejme ještě jedno nové tvrzení.

Věta 2.3 (Vlastnosti relativní četnosti)

- (i) $E \bar{X}_n = p$, $\text{var } \bar{X}_n = \frac{p(1-p)}{n}$;

* Angl. *Bernoulli distribution*. † Angl. *empirical frequency*

- (ii) $\bar{X}_n \xrightarrow{P} p$ pro $n \rightarrow \infty$;
 (iii) $\sqrt{n}(\bar{X}_n - p) \xrightarrow{d} N(0, p(1-p))$ pro $n \rightarrow \infty$;
 (iv) $n\bar{X}_n \sim \text{Bi}(n, p)$.

Důkaz. (i) až (iii) plyne přímo z věty 2.2 s využitím toho, že pro alternativní náhodnou veličinu platí $E X_i = p$ a $\text{var } X_i = p(1-p)$. (iv) plyne z rovnosti $n\bar{X}_n = \sum_{i=1}^n X_i$ a z reprezentace binomického rozdělení jako součtu nezávislých stejně rozdělených alternativních rozdělení. \square

Podle bodu (ii) můžeme pravděpodobnost p zjistit s libovolnou přesností pomocí relativní četnosti, stačí jen mít dostatek pozorování výskytu tohoto jevu.

2.2.3. VLASTNOSTI VÝBĚROVÉHO ROZPTYLU

Nejprve uvažujme obecný model $\mathcal{F} = \mathcal{L}^2$. Označme opět $\mu = E X_i$ a $\sigma^2 = \text{var } X_i$. Výběrový rozptyl lze přepsat do různých podob, které se k určitým účelům hodí lépe než původní definice.

Věta 2.4 Pro $n \geq 2$

(i)

$$S_n^2 = \frac{n}{n-1} \left(\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 \right). \quad (2.1)$$

(ii) Nechť $\mathbf{1}_n$ je sloupcový vektor n jedniček. Označme $\mathbb{A} = \mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$ (matice $n \times n$).

Pak

$$S_n^2 = \frac{1}{n-1} \mathbf{X}^T \mathbb{A} \mathbf{X} = \frac{1}{n-1} \mathbf{Y}^T \mathbb{A} \mathbf{Y}, \quad (2.2)$$

kde $\mathbf{Y} = \mathbf{X} - c \mathbf{1}_n$ pro nějaké $c \in \mathbb{R}$.

Důkaz. Část (i):

$$\begin{aligned} \frac{n-1}{n} S_n^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n} \sum_{i=1}^n (X_i^2 - 2X_i \bar{X}_n + \bar{X}_n^2) = \frac{1}{n} \sum_{i=1}^n X_i^2 - \frac{2}{n} \sum_{i=1}^n X_i \bar{X}_n + \bar{X}_n^2 \\ &= \frac{1}{n} \sum_{i=1}^n X_i^2 - 2\bar{X}_n^2 + \bar{X}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2. \end{aligned}$$

Část (ii):

$$\begin{aligned} \mathbf{X}^T \mathbb{A} \mathbf{X} &= \mathbf{X}^T \left(\mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right) \mathbf{X} = \mathbf{X}^T \mathbf{X} - \frac{1}{n} \mathbf{X}^T \mathbf{1}_n \mathbf{1}_n^T \mathbf{X} \\ &= \sum_{i=1}^n X_i^2 - \frac{1}{n} \left(\sum_{i=1}^n X_i \right)^2 = \sum_{i=1}^n X_i^2 - n\bar{X}_n^2 = (n-1)S_n^2. \end{aligned}$$

Poslední část tvrzení pak plyne z toho, že

$$\mathbf{1}_n^\top \mathbb{A} = \mathbf{0} = \mathbb{A} \mathbf{1}_n.$$

□

Poznámka. Vzorec (2.1) se používá mj. pro numerický výpočet S_n^2 . Vzorec (2.2) přepisuje S_n^2 v podobě kvadratické formy a ukazuje, že S_n^2 je invariantní vůči posunutí pozorování X_i o libovolnou konstantu c .

Povšimněte si, že matice \mathbb{A} je idempotentní, tj. $\mathbb{A}\mathbb{A} = \mathbb{A}$. To využijeme později při určování rozdělení S_n^2 (viz věta 2.8 níže).

U kvadratických forem máme k dispozici šikovný vzorec pro výpočet střední hodnoty.

Lemma 2.5 Nechť Z je náhodný vektor délky n se střední hodnotou μ a konečnou rozptylovou maticí Σ . Nechť \mathbb{B} je libovolná matice $n \times n$. Pak platí

$$E Z^\top \mathbb{B} Z = \mu^\top \mathbb{B} \mu + \text{tr}(\mathbb{B} \Sigma).$$

Důkaz.

$$\begin{aligned} E Z^\top \mathbb{B} Z &= E \text{tr}(Z^\top \mathbb{B} Z) = E \text{tr}(\mathbb{B} Z Z^\top) = \text{tr}(\mathbb{B} E Z Z^\top) = \text{tr}(\mathbb{B}(\mu \mu^\top + \Sigma)) \\ &= \text{tr}(\mathbb{B} \mu \mu^\top) + \text{tr}(\mathbb{B} \Sigma) = \mu^\top \mathbb{B} \mu + \text{tr}(\mathbb{B} \Sigma), \end{aligned}$$

kde jsme využili toho, že

$$\Sigma = E(Z - \mu)(Z - \mu)^\top = E Z Z^\top - \mu \mu^\top.$$

□

Věta 2.6 (Vlastnosti výběrového rozptylu)

(i) $S_n^2 \xrightarrow[n \rightarrow \infty]{P} \sigma^2$.

(ii) $E S_n^2 = \sigma^2$.

(iii) Jestliže $\mathcal{F} = \mathcal{L}^4$ (existuje konečný čtvrtý moment X_i), pak

$$\sqrt{n} (S_n^2 - \sigma^2) \xrightarrow[n \rightarrow \infty]{d} N(0, \sigma^4(\gamma_4 - 1)),$$

kde $\gamma_4 = \frac{E(X_i - \mu)^4}{\sigma^4}$ je tzv. špičatost* rozdělení X_i .

* Angl. *kurtosis*

(iv) Jestliže $\mathcal{F} = \mathcal{L}^4$, pak

$$\sqrt{n} \left[\begin{pmatrix} \bar{X}_n \\ S_n^2 \end{pmatrix} - \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix} \right] \xrightarrow[n \rightarrow \infty]{d} N_2(\mathbf{0}, \Sigma),$$

kde $\Sigma = \begin{pmatrix} \sigma^2 & \sigma^3 \gamma_3 \\ \sigma^3 \gamma_3 & \sigma^4 (\gamma_4 - 1) \end{pmatrix}$ a $\gamma_3 = \frac{E(X_i - \mu)^3}{\sigma^3}$ je tzv. šikmost* rozdělení X_i .

Důkaz. Část (i): Dle věty 2.4(i) můžeme psát

$$S_n^2 = \frac{n}{n-1} \left(\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 \right).$$

Jelikož $\frac{n}{n-1} \xrightarrow[n \rightarrow \infty]{} 1$, tak stačí dokázat, že

$$\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 \xrightarrow[n \rightarrow \infty]{P} \sigma^2.$$

Ze zákona velkých čísel (tvrzení 1.5) platí

$$\left(\bar{X}_n, \frac{1}{n} \sum_{i=1}^n X_i^2 \right)^T \xrightarrow[n \rightarrow \infty]{P} \left(E X_i, E X_i^2 \right)^T.$$

Nyní funkce $g(y_1, y_2) = y_2 - y_1^2$ je spojitá na \mathbb{R}^2 , tedy je spojitá i v daném (neznámém bodě) $(E X_i, E X_i^2)$, který je nosičem limitního rozdělení. Tedy můžeme použít větu o spojitě transformaci (tvrzení 1.2(ii)) a dostáváme

$$\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 \xrightarrow[n \rightarrow \infty]{P} E X_i^2 - (E X_i)^2 = \text{var } X_i = \sigma^2,$$

což jsme měli ověřit.

Část (ii): Položme $Y = X - \mu \mathbf{1}_n$ a všimněme si, že $E Y = \mathbf{0}$. Dle věty 2.4(ii) a lemmatu 2.5 můžeme počítat

$$(n-1) E S_n^2 = E Y^T A Y = E Y^T A E Y + \text{tr}(A \sigma^2 \mathbb{1}_n) = \mathbf{0} + (n-1) \sigma^2,$$

neboť

$$\text{tr}(A \sigma^2 \mathbb{1}_n) = \sigma^2 \left(\text{tr}(\mathbb{1}_n) - \frac{1}{n} \text{tr}(\mathbf{1}_n \mathbf{1}_n^T) \right) = \sigma^2 (n-1).$$

Část (iii): Nejdříve si přepíšeme výběrový rozptyl jako

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu)^2 - \frac{n}{n-1} (\bar{X}_n - \mu)^2.$$

* Angl. *skewness*

A tedy

$$\sqrt{n}(S_n^2 - \sigma^2) = \frac{\sqrt{n}}{n-1} \sum_{i=1}^n [(X_i - \mu)^2 - \sigma^2] + \frac{\sqrt{n}}{n-1} \sigma^2 - \frac{n}{n-1} \sqrt{n} (\bar{X}_n - \mu)^2 \stackrel{\text{ozn.}}{=} A_n + B_n + C_n,$$

kde A_n , B_n a C_n postupně značí jednotlivé sčítance na pravé straně rovnosti. Zřejmě

$$B_n = \frac{\sqrt{n}}{n-1} \sigma^2 \xrightarrow{n \rightarrow \infty} 0.$$

Dále

$$C_n = \frac{n}{n-1} \sqrt{n} (\bar{X}_n - \mu)^2 = \frac{n}{n-1} \sqrt{n} (\bar{X}_n - \mu) (\bar{X}_n - \mu) \xrightarrow[n \rightarrow \infty]{P} 0,$$

kde jsme využili toho, že $\frac{n}{n-1} \xrightarrow{n \rightarrow \infty} 1$, $\sqrt{n} (\bar{X}_n - \mu) \xrightarrow[n \rightarrow \infty]{d} N(0, \sigma^2)$, $\bar{X}_n - \mu \xrightarrow[n \rightarrow \infty]{P} 0$ a Cramérov-Sluckého věty (tvrzení 1.3).

Stačí se tedy zabývat členem A_n . Pro $i \in \{1, \dots, n\}$ si označme $Y_i = (X_i - \mu)^2$. Potom s využitím centrální limitní věty na náhodné veličiny Y_i (tvrzení 1.6) a Cramérov-Sluckého věty (tvrzení 1.3)

$$\begin{aligned} A_n &= \frac{\sqrt{n}}{n-1} \sum_{i=1}^n [(X_i - \mu)^2 - \sigma^2] = \frac{n}{n-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n [(X_i - \mu)^2 - \sigma^2] \\ &= \frac{n}{n-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n [Y_i - E Y_i] \xrightarrow[n \rightarrow \infty]{d} N(0, \text{var}(Y_i)). \end{aligned}$$

Zbývá už jen dopočítat

$$\text{var}(Y_i) = \text{var}((X_i - \mu)^2) = E(X_i - \mu)^4 - (\sigma^2)^2 = \sigma^4 [E\left(\frac{X_i - \mu}{\sigma}\right)^4 - 1] = \sigma^4 [\gamma_4 - 1].$$

□

Poznámka.

- Věta 2.6(iii) říká, že variabilita výběrového rozptylu asymptoticky závisí na špičatosti pozorování.
- Věta 2.6(iv) říká, že výběrový průměr a výběrový rozptyl mají asymptoticky sdružené normální rozdělení. Jejich kovariance asymptoticky závisí na šikmosti pozorování. Je-li šikmost nulová, výběrový průměr a výběrový rozptyl jsou asymptoticky nezávislé.

Poznámka. Alternativně se věta 2.6(ii) (tj. nestrannost výběrového rozptylu) dá ukázat přímočarým výpočtem.

$$\begin{aligned} E S_n^2 &= \frac{1}{n-1} \left(\sum_{i=1}^n E X_i^2 - n E \bar{X}_n^2 \right) = \frac{1}{n-1} \left(n E X_1^2 - n \text{var}(\bar{X}_n) - n (E \bar{X}_n)^2 \right) \\ &= \frac{1}{n-1} \left(n(\sigma^2 + \mu^2) - n \frac{\sigma^2}{n} - n \mu^2 \right) = \frac{1}{n-1} (n\sigma^2 - \sigma^2) = \sigma^2, \end{aligned}$$

kde jsme využili toho, že $E X_1^2 = \text{var}(X_1) + (E X_1)^2$ a podobně také $E (\bar{X}_n)^2 = \text{var}(\bar{X}_n) + (E \bar{X}_n)^2$.

Cvičení. Dokažte, že pokud X_i nabývají pouze hodnot 0 nebo 1, pak $S_n^2 = \frac{n}{n-1} \bar{X}_n (1 - \bar{X}_n)$. *Návod: Využijte toho, že v tomto případě $X_i^2 = X_i$.*

Nyní přidáme **předpoklad normálního rozdělení**, tj. budeme pracovat v menším modelu $\mathcal{F} = \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}^+\}$. Pracujeme tedy s náhodným výběrem $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$, kde X_i jsou nezávislé s rozdělením $N(\mu, \sigma^2)$. Díky jejich nezávislosti platí $\mathbf{X} \sim N_n(\mu \mathbf{1}_n, \sigma^2 \mathbb{I}_n)$.

Nejprve uvedeme dva výsledky, které platí pro libovolné normálně rozdělené náhodné vektory.

Lemma 2.7 Nechť $\mathbf{X} \sim N_n(\boldsymbol{\mu}, \Sigma)$ a \mathbb{A} je pozitivně semidefinitní matice typu $n \times n$.

- (i) Nechť \mathbb{B} je libovolná matice typu $m \times n$ splňující rovnost $\mathbb{B}\Sigma\mathbb{A} = \mathbb{0}_{m \times n}$. Pak náhodná veličina $\mathbf{X}^T \mathbb{A} \mathbf{X}$ a náhodný vektor $\mathbb{B} \mathbf{X}$ jsou nezávislé.
- (ii) Nechť \mathbb{B} je libovolná pozitivně semidefinitní matice typu $n \times n$ splňující rovnost $\mathbb{B}\Sigma\mathbb{A} = \mathbb{0}_{n \times n}$. Pak jsou náhodné veličiny $\mathbf{X}^T \mathbb{A} \mathbf{X}$ a $\mathbf{X}^T \mathbb{B} \mathbf{X}$ nezávislé.

Důkaz. Část (i). Jelikož matice \mathbb{A} je pozitivně semidefinitní, tudíž dle spektrálního rozkladu existuje ortonormální matice \mathbb{U} taková, že

$$\mathbb{A} = \mathbb{U} \mathbb{D} \mathbb{U}^T$$

kde $\mathbb{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$ je diagonální s vlastními čísly matice \mathbb{A} (která jsou nezáporná).

Dále z předpokladu věty máme

$$\mathbb{0}_{m \times n} = \mathbb{B}\Sigma\mathbb{A} = \mathbb{B}\Sigma\mathbb{U}\mathbb{D}\mathbb{U}^T.$$

Jako $\mathbb{D}^{-1/2}$ označme nyní diagonální matici, která má na i -tém prvku diagonály $\frac{1}{\sqrt{\lambda_i}}$ v případě, že λ_i je nenulové a nulu jinak. Pak vynásobením výše uvedené rovnosti zprava maticí $\mathbb{U}\mathbb{D}^{-1/2}$ dostáváme

$$\mathbb{0}_{m \times n} = \mathbb{B}\Sigma\mathbb{U}\mathbb{D}^{1/2}.$$

Tedy náhodné vektory $\mathbb{B} \mathbf{X}$ a $\mathbb{D}^{1/2} \mathbb{U}^T \mathbf{X}$ jsou nekorelované, neboť

$$\text{cov}(\mathbb{B} \mathbf{X}, \mathbb{D}^{1/2} \mathbb{U}^T \mathbf{X}) = \mathbb{B}\Sigma\mathbb{U}\mathbb{D}^{1/2} = \mathbb{0}_{m \times n}.$$

Z definice mnohorozměrného normálního rozdělení plyne, že tyto náhodné vektory mají sdruženě normální rozdělení, neboť můžeme psát

$$\begin{pmatrix} \mathbb{B} \mathbf{X} \\ \mathbb{D}^{1/2} \mathbb{U}^T \mathbf{X} \end{pmatrix} = \begin{pmatrix} \mathbb{B} \\ \mathbb{D}^{1/2} \mathbb{U}^T \end{pmatrix} \mathbf{X}.$$

Sdružená normalita a nekorelovanost pak implikuje nezávislost náhodných vektorů $\mathbb{B} \mathbf{X}$ a $\mathbb{D}^{1/2} \mathbb{U}^T \mathbf{X}$ (věta P.6.2(ii)). Tudíž také $\mathbb{B} \mathbf{X}$ a $\mathbf{X}^T \mathbb{U} \mathbb{D}^{1/2} \mathbb{D}^{1/2} \mathbb{U}^T \mathbf{X} = \mathbf{X}^T \mathbb{A} \mathbf{X}$ jsou nezávislé.

Část (ii). Podobně jako výše ze spektrálního rozkladu

$$\mathbb{A} = \mathbb{U}_A \mathbb{D}_A \mathbb{U}_A^T \quad \text{a} \quad \mathbb{B} = \mathbb{U}_B \mathbb{D}_B \mathbb{U}_B^T,$$

kde $\mathbb{U}_A, \mathbb{U}_B$ jsou ortonormální matice a $\mathbb{D}_A, \mathbb{D}_B$ jsou diagonální matice s nezápornými prvky na diagonálách.

Dále z předpokladu

$$\mathbb{0}_{n \times n} = \mathbb{B} \Sigma \mathbb{A} = \mathbb{U}_B \mathbb{D}_B \mathbb{U}_B^T \Sigma \mathbb{U}_A \mathbb{D}_A \mathbb{U}_A^T$$

Nechť $\mathbb{D}_A^{-1/2}$ a $\mathbb{D}_B^{-1/2}$ jsou jako matice $\mathbb{D}^{-1/2}$ výše. Pak vynásobením výše uvedené rovnosti zprava maticí $\mathbb{U}_A \mathbb{D}_A^{-1/2}$ a zleva maticí $\mathbb{D}_B^{-1/2} \mathbb{U}_B^T$ dostáváme

$$\mathbb{0}_{n \times n} = \mathbb{D}_B^{1/2} \mathbb{U}_B^T \Sigma \mathbb{U}_A \mathbb{D}_A^{1/2}.$$

Tedy podobně jako v části (i) dostáváme, že náhodné vektory $\mathbb{D}_B^{1/2} \mathbb{U}_B^T \mathbf{X}$ a $\mathbb{D}_A^{1/2} \mathbb{U}_A^T \mathbf{X}$ jsou nezávislé a tudíž také jejich skalární součiny

$$\mathbf{X}^T \mathbb{U}_B \mathbb{D}_B^{1/2} \mathbb{D}_B^{1/2} \mathbb{U}_B^T \mathbf{X} = \mathbf{X}^T \mathbb{B} \mathbf{X}$$

a

$$\mathbf{X}^T \mathbb{U}_A \mathbb{D}_A^{1/2} \mathbb{D}_A^{1/2} \mathbb{U}_A^T \mathbf{X} = \mathbf{X}^T \mathbb{A} \mathbf{X}.$$

jsou nezávislé. □

Věta 2.8 (Vlastnosti výběrového rozptylu za normality) Nechť $X_i \sim N(\mu, \sigma^2)$, $i = 1, \dots, n$ jsou nezávislé. Pak platí

(i)

$$\frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{n-1}^2. \quad (2.3)$$

(ii) \bar{X}_n a S_n^2 jsou nezávislé náhodné veličiny.

Důkaz. Část (i). Dle věty 2.4 můžeme psát

$$\frac{(n-1)S_n^2}{\sigma^2} = \mathbf{Y}^T \mathbb{A} \mathbf{Y},$$

kde

$$\mathbf{Y} = \left(\frac{X_1 - \mu}{\sigma}, \dots, \frac{X_n - \mu}{\sigma} \right)^T \sim \mathbb{N}_n(\mathbf{0}, \mathbb{I}_n)$$

a $\mathbb{A} = \mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$. Jelikož matice \mathbb{A} je idempotentní s hodnotí $n-1$, tak tvrzení plyne z lemmatu A.4 (kde $\Sigma = \mathbb{I}_n$).

* Viz definice A.1.

Část (ii) Všimněme si, že můžeme psát

$$\bar{X}_n = \frac{1}{n} \mathbb{B} \mathbf{X}, \quad S_n^2 = \frac{1}{n-1} \mathbf{X}^\top \mathbb{A} \mathbf{X},$$

kde $\mathbb{B} = \mathbf{1}_n^\top$ a $\mathbb{A} = \mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$. Dále $\mathbf{X} \sim \mathbb{N}_n(\mu \mathbf{1}_n, \sigma^2 \mathbb{I}_n)$ a tedy tvrzení plyne z lemmatu 2.7(i), neboť

$$\mathbb{B} \Sigma \mathbb{A} = \mathbf{1}_n^\top \sigma^2 \mathbb{I}_n (\mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) = \sigma^2 (\mathbf{1}_n^\top - \frac{1}{n} n \mathbf{1}_n^\top) = \mathbf{0}_n^\top.$$

□

Poznámka. Z definice χ^2 rozdělení víme, že náhodná veličina s rozdělením χ_{n-1}^2 má rozdělení dané pomocí $\sum_{i=1}^{n-1} Y_i^2$, kde Y_1, \dots, Y_{n-1} jsou nezávislé stejně rozdělené náhodné veličiny s rozdělením $N(0, 1)$. Z centrální limitní věty a z (2.3) pak plyne, že

$$\frac{\frac{(n-1)S_n^2}{\sigma^2} - (n-1)}{\sqrt{n-1}} \xrightarrow[n \rightarrow \infty]{d} N(0, 2)$$

a tudíž

$$\sqrt{\frac{n-1}{n}} \sqrt{n} (S_n^2 - \sigma^2) \stackrel{\text{as.}}{\sim} N(0, 2\sigma^4).$$

Uvědomíme-li si, že špičatost normálního rozdělení je 3, vidíme, že tvrzení (i) z věty 2.8 je v souladu s asymptotickým výsledkem věty 2.6(iii). Věta 2.8(i) udává přesné rozdělení S_n^2 pro normální data, zatímco věta 2.6(iii) udává asymptotické rozdělení S_n^2 pro libovolná data s konečným čtvrtým momentem.

Poznámka. Tvrzení věty 2.8(i) se dá pamatovat následovně. Všimněme si, že

$$\frac{(n-1)S_n^2}{\sigma^2} = \sum_{i=1}^n \left(\frac{X_i - \bar{X}_n}{\sigma} \right)^2.$$

Kdybychom ve výrazu na pravé straně předchozí rovnosti použili skutečnou střední hodnotu μ místo \bar{X}_n , tak bychom měli $\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 \sim \chi_n^2$. Odhadnutím střední hodnoty μ pomocí \bar{X}_n pak jakoby ztrácíme jeden stupeň volnosti (protože jsme odhadovali jeden parametr).

Poznámka. Věta 2.8(ii) říká, že jsou-li data normální, \bar{X}_n a S_n^2 jsou nezávislé pro každé konečné $n > 1$.

Věta 2.9 (limitní věta o T_n) Nechť X_1, \dots, X_n je náhodný výběr z libovolného rozdělení se střední hodnotou μ a s konečným nenulovým rozptylem σ^2 . Pak

$$T_n = \frac{\sqrt{n} (\bar{X}_n - \mu)}{S_n} \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Důkaz. Náhodnou veličinu T_n si můžeme přepsat do tvaru

$$T_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \frac{\sigma}{S_n}.$$

Z centrální limitní věty (tvrzení 1.6, pro $k = 1$) máme, že

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Dále z $S_n^2 \xrightarrow[n \rightarrow \infty]{P} \sigma^2$ (věta 2.6(i)) a z věty o spojitě transformaci (tvrzení 1.2(ii)) pro $g(y) = \sigma/\sqrt{y}$ plyne, že

$$\frac{\sigma}{S_n} \xrightarrow[n \rightarrow \infty]{P} 1.$$

Tvrzení pak plyne z Cramérový-Sluckého věty (tvrzení 1.3). □

Nyní opět přidáme předpoklad normálního rozdělení.

Věta 2.10 Necht' X_1, \dots, X_n je náhodný výběr z rozdělení $N(\mu, \sigma^2)$. Pak

$$T_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} \sim t_{n-1}.*$$

Důkaz. Náhodnou veličinu T_n si můžeme přepsat do tvaru

$$T_n = \frac{\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma}}{\sqrt{\frac{(n-1)S_n^2}{\sigma^2} / (n-1)}}. \quad (2.4)$$

Z poznámky za větou 2.2 víme, že $\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \sim N(0, 1)$. Dále $\frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{n-1}^2$ (věta 2.8(i)), přičemž čitatel a jmenovatel ve zlomku (2.4) jsou nezávislé (věta 2.8(ii)). Tvrzení pak plyne z reprezentace t -rozdělení (Definice A.2). □

Poznámka. Věta 2.10 udává přesné rozdělení náhodné veličiny T_n pro normální data, zatímco věta 2.9 udává asymptotické rozdělení téže veličiny pro libovolná data s konečným rozptylem. Uvědomte si, že pro $n \rightarrow \infty$ rozdělení t_{n-1} konverguje v distribuci k rozdělení $N(0, 1)$.

Nyní budeme uvažovat dva nezávislé výběry ze dvou různých normálních rozdělení.

Definice 2.5 (o F -rozdělení) Necht' $X \sim \chi_n^2$ a $Y \sim \chi_m^2$ jsou nezávislé. Pak rozdělení náhodné veličiny

$$Z = \frac{X/n}{Y/m}$$

se nazývá [Fisherovo-Snedecorovo] F -rozdělení s n a m stupni volnosti, značíme $F_{n,m}$.

* Viz definice A.2.

Věta 2.11 (věta o F statistic) Nechť X_1, \dots, X_n je náhodný výběr z normálního rozdělení $N(\mu_X, \sigma_X^2)$ a Y_1, \dots, Y_m je náhodný výběr z normálního rozdělení $N(\mu_Y, \sigma_Y^2)$. Nechť jsou náhodné vektory $(X_1, \dots, X_n)^T$ a $(Y_1, \dots, Y_m)^T$ nezávislé. Označme výběrové průměry obou výběrů \bar{X}_n a \bar{Y}_m a výběrové rozptyly

$$S_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \quad \text{a} \quad S_Y^2 = \frac{1}{m-1} \sum_{j=1}^m (Y_j - \bar{Y}_m)^2.$$

Pak platí

$$\frac{S_X^2/\sigma_X^2}{S_Y^2/\sigma_Y^2} \sim F_{n-1, m-1}.$$

Důkaz. Statistiku si můžeme přepsat jako

$$\frac{S_X^2/\sigma_X^2}{S_Y^2/\sigma_Y^2} = \frac{\frac{(n-1)S_X^2}{\sigma_X^2}/(n-1)}{\frac{(m-1)S_Y^2}{\sigma_Y^2}/(m-1)}.$$

Dále $\frac{(n-1)S_X^2}{\sigma_X^2} \sim \chi_{n-1}^2$ a $\frac{(m-1)S_Y^2}{\sigma_Y^2} \sim \chi_{m-1}^2$ (věta 2.8(ii)), přičemž tyto náhodné veličiny jsou nezávislé. Tvrzení pak plyne z definice F-rozdělení (definice 2.5). \square

2.3. USPOŘÁDANÝ NÁHODNÝ VÝBĚR

Mějme náhodný výběr X_1, \dots, X_n z jednorozměrného spojitého rozdělení s distribuční funkcí F a hustotou f vzhledem k Lebesgueově míře. Nechť $n \geq 2$. Jelikož X_1, \dots, X_n jsou nezávislé a mají spojitě rozdělení, tak

$$P(X_i = X_j \text{ pro nějaká } i, j \in \{1, \dots, n\}) = 0.$$

Definice 2.6 (Uspořádaný náhodný výběr a pořadí)

- (i) Seřadíme-li všechny náhodné veličiny* X_1, \dots, X_n od nejmenší do největší, získáme *uspořádaný náhodný výběr*†

$$X_{(1)} < X_{(2)} < \dots < X_{(n-1)} < X_{(n)}.$$

Symbolem $X_{(k)}$ rozumíme k -tou nejmenší hodnotu mezi pozorováními X_1, \dots, X_n ; nazýváme ji k -tá *pořádková statistika*‡.

- (ii) *Pořadím*§ náhodné veličiny X_i ve výběru X_1, \dots, X_n rozumíme přirozené číslo $R_i \in \{1, \dots, n\}$ takové, že $X_i = X_{(R_i)}$.

* Seřazujeme vlastně realizace náhodných veličin, více poznámka pod definicí. † Angl. *ordered random sample* ‡ Angl. *order statistic* § Angl. *rank*

Celý uspořádaný výběr budeme značit $\mathbf{X}_{(\cdot)}$, tj.

$$\mathbf{X}_{(\cdot)} = (X_{(1)}, \dots, X_{(n)})^\top.$$

Poznámka.

1. Jelikož náhodné veličiny X_1, \dots, X_n jsou zobrazení z Ω do \mathbb{R} , tak pro každé $\omega \in \Omega$ vlastně seřazujeme jejich realizace $X_1(\omega), \dots, X_n(\omega)$. Tedy např. $X_{(1)}$ je opět zobrazení z Ω do \mathbb{R} definované pro $\omega \in \Omega$ jako $X_{(1)}(\omega) = \min\{X_1(\omega), \dots, X_n(\omega)\}$.
2. Hodnoty X_1, \dots, X_n lze jednoznačně určit z n -tice pořádkových statistik a n -tice pořadí.
3. První pořádková statistika je minimum, n -tá pořádková statistika je maximum všech veličin náhodného výběru.
4. Platí $R_i = \sum_{j=1}^n \mathbb{1}\{X_i \geq X_j\} = 1 + \sum_{j=1}^n \mathbb{1}\{X_i > X_j\}$.
5. Pořádkové statistiky a pořadí jsou náhodné veličiny a též statistiky ve smyslu definice 2.3.

Označme symbolem \mathcal{P}_n množinu všech permutací posloupnosti $(1, \dots, n)$. Tato množina má $n!$ prvků.

Věta 2.12 Sdružená hustota náhodného vektoru $\mathbf{X}_{(\cdot)} = (X_{(1)}, \dots, X_{(n)})^\top$ vzhledem k Lebesgueově míře jest

$$p(y_1, \dots, y_n) = \begin{cases} n! f(y_1) f(y_2) \cdots f(y_n) & \text{pokud } y_1 < \cdots < y_n, \\ 0 & \text{jinak.} \end{cases}$$

Důkaz. Víme, že náhodný vektor $\mathbf{X}_{(\cdot)}$ má hustotu p , právě když pro každou borelovskou množinu $B \in \mathcal{B}_0^n$ platí

$$P(\mathbf{X}_{(\cdot)} \in B) = \int \cdots \int \mathbb{1}_B(\mathbf{y}) p(\mathbf{y}) d\mathbf{y}.$$

Přednášková verze zbytku důkazu.

Vezměme si tedy $B \in \mathcal{B}_0^n$ a pro jednoduchost uvažujme $n = 2$. Potom

$$\begin{aligned} P(\mathbf{X}_{(\cdot)} \in B) &= P((X_{(1)}, X_{(2)})^\top \in B) \\ &= P((X_{(1)}, X_{(2)})^\top \in B, X_1 < X_2) + P((X_{(1)}, X_{(2)})^\top \in B, X_2 < X_1) \\ &= P((X_1, X_2)^\top \in B, X_1 < X_2) + P((X_2, X_1)^\top \in B, X_2 < X_1). \end{aligned}$$

Nyní si všimněme, že nezávislost a stejná rozdělenost náhodných vektorů implikuje, že rozdělení náhodného vektoru $(X_1, X_2)^\top$ je stejné jako rozdělení náhodného vektoru $(X_2, X_1)^\top$. Tudíž

$$P((X_1, X_2)^\top \in B, X_1 < X_2) = P((X_2, X_1)^\top \in B, X_2 < X_1)$$

a tedy dostáváme

$$\begin{aligned} P(\mathbf{X}_{(\cdot)} \in B) &= 2 P((X_1, X_2)^\top \in B, X_1 < X_2) \\ &= 2 \iint \mathbb{1}_B(y_1, y_2) f(y_1) f(y_2) \mathbb{1}_{\{y_1 < y_2\}} dy_1 dy_2. \end{aligned}$$

Pro obecné n pak plyne důkaz z toho, že obdobně jako výše se dokáže, že

$$P(\mathbf{X}_{(\cdot)} \in B) = n! P(\mathbf{X} \in B, X_1 < X_2 < \dots < X_n).$$

Skriptová verze zbytku důkazu.

Vezměme si tedy $B \in \mathcal{B}_0^n$ a označme si vektor pořadí $\mathbf{R} = (R_1, \dots, R_n)^\top$ a jednu z permutací množiny \mathcal{P}_n jako $\mathbf{r} = (r_1, \dots, r_n)^\top$. Je dobré si uvědomit, že \mathbf{R} závisí na \mathbf{X} . Proto tam, kde to bude vhodné, tak budeme psát $\mathbf{R}(\mathbf{X})$.

Nyní můžeme počítat

$$\begin{aligned} P(\mathbf{X}_{(\cdot)} \in B) &= \sum_{\mathbf{r} \in \mathcal{P}_n} P(\mathbf{X}_{(\cdot)} \in B, \mathbf{R}(\mathbf{X}) = \mathbf{r}) \\ &= \sum_{\mathbf{r} \in \mathcal{P}_n} \int \dots \int \mathbb{1}_B(\mathbf{x}_{(\cdot)}) \mathbb{1}\{\mathbf{R}(\mathbf{x}) = \mathbf{r}\} f(x_1) \dots f(x_n) dx_1 \dots dx_n \\ &= \sum_{\mathbf{r} \in \mathcal{P}_n} \int \dots \int \mathbb{1}_B(\mathbf{y}) \mathbb{1}\{\mathbf{R}(y_{r_1}, \dots, y_{r_n}) = \mathbf{r}\} f(y_{r_1}) \dots f(y_{r_n}) dy_1 \dots dy_n \\ &= \int \dots \int \mathbb{1}_B(\mathbf{y}) \mathbb{1}\{y_1 < \dots < y_n\} \sum_{\mathbf{r} \in \mathcal{P}_n} f(y_{r_1}) \dots f(y_{r_n}) dy_1 \dots dy_n \\ &= \int \dots \int \mathbb{1}_B(\mathbf{y}) \mathbb{1}\{y_1 < \dots < y_n\} n! f(y_1) \dots f(y_n) dy_1 \dots dy_n, \end{aligned}$$

kde jsme přeznačili $\mathbf{y} = \mathbf{x}_{(\cdot)}$ (tj. $x_i = y_{r_i}$, $i = 1, \dots, n$), z čehož plyne, že $y_1 < \dots < y_n$. Tudiž také hodnoty y_{r_1}, \dots, y_{r_n} mají pořadí \mathbf{r} a tedy indikátor $\mathbb{1}\{\mathbf{R}(y_{r_1}, \dots, y_{r_n}) = \mathbf{r}\}$ je splněn a mohli jsme ho nahradit indikátorem $\mathbb{1}\{y_1 < \dots < y_n\}$. \square

Poznámka. Náhodné veličiny $X_{(1)}, \dots, X_{(n)}$ nejsou nezávislé. Podobně ani náhodné veličiny udávající pořadí R_1, \dots, R_n nejsou nezávislé.

Věta 2.13 Distribuční funkce k -té pořádkové statistiky jest

$$\begin{aligned} F_{(k)}(x) &= P(X_{(k)} \leq x) = \sum_{j=k}^n \binom{n}{j} F^j(x) (1 - F(x))^{n-j} \\ &= \frac{1}{B(k, n - k + 1)} \int_0^{F(x)} t^{k-1} (1 - t)^{n-k} dt, \end{aligned}$$

kde $B(\cdot, \cdot)$ značí Beta funkci (viz Appendix A.5).

2. Náhodný výběr

Důkaz. První rovnost: Označme si $Z_i = \mathbb{1}\{X_i \leq x\}$. Potom $Y_n = \sum_{i=1}^n Z_i$ udává počet veličin, které jsou menší než x . Navíc $Y_n \sim \text{Bi}(n, F(x))$. Tudíž

$$P(X_{(k)} \leq x) = P(Y_n \geq k) = \sum_{j=k}^n P(Y_n = j) = \sum_{j=k}^n \binom{n}{j} F^j(x) (1 - F(x))^{n-j}.$$

Druhá rovnost: Budeme postupovat zpětnou indukcí.

Nechť $k = n$, potom

$$\frac{1}{B(n, 1)} \int_0^{F(x)} t^{n-1} dt = \frac{1}{B(n, 1)} \frac{1}{n} F^n(x) = \binom{n}{n} F^n(x),$$

kde jsme využili toho, že

$$\frac{1}{n B(n, 1)} = \frac{\Gamma(n+1)}{n \Gamma(n) \Gamma(1)} = 1 = \binom{n}{n}.$$

Nyní provedeme *indukční krok* ($k \rightarrow k-1$). Předpokládejme, že pro dané k platí

$$\sum_{j=k}^n \binom{n}{j} F^j(x) (1 - F(x))^{n-j} = \frac{1}{B(k, n-k+1)} \int_0^{F(x)} t^{k-1} (1-t)^{n-k} dt$$

a chceme ukázat, že

$$\sum_{j=k-1}^n \binom{n}{j} F^j(x) (1 - F(x))^{n-j} = \frac{1}{B(k-1, n-k+2)} \int_0^{F(x)} t^{k-2} (1-t)^{n-k+1} dt. \quad (2.5)$$

Počítejme nyní pomocí metody per partes integrál na pravé straně předcházející rovnosti, tj.

$$\begin{aligned} \int_0^{F(x)} t^{k-2} (1-t)^{n-k+1} dt &= \left[\frac{1}{k-1} t^{k-1} (1-t)^{n-k+1} \right]_0^{F(x)} + \frac{n-k+1}{k-1} \int_0^{F(x)} t^{k-1} (1-t)^{n-k} dt \\ &= \frac{1}{k-1} F^{k-1}(x) (1-F(x))^{n-k+1} + \frac{n-k+1}{k-1} \int_0^{F(x)} t^{k-1} (1-t)^{n-k} dt. \end{aligned}$$

Pravá strana (2.5) se tedy rovná

$$\frac{1}{B(k-1, n-k+2)(k-1)} \left(F^{k-1}(x) (1-F(x))^{n-k+1} + (n-k+1) \int_0^{F(x)} t^{k-1} (1-t)^{n-k} dt \right).$$

Nyní si všimněme, že

$$\frac{1}{B(k-1, n-k+2)(k-1)} = \frac{\Gamma(n+1)}{\Gamma(k-1)\Gamma(n-k+2)(k-1)} = \frac{n!}{(k-1)!(n-k+1)!} = \binom{n}{k-1}$$

a dále

$$\frac{n-k+1}{B(k-1, n-k+2)(k-1)} = \frac{\Gamma(n+1)}{\Gamma(k)\Gamma(n-k+1)} = \frac{1}{B(k, n-k+1)}.$$

Odtud již s využitím indukčního předpokladu dostáváme pro pravou stranu (2.5), že

$$\begin{aligned} & \frac{1}{B(k-1, n-k+2)} \int_0^{F(x)} t^{k-2}(1-t)^{n-k+1} dt \\ &= \binom{n}{k-1} F^{k-1}(x)(1-F(x))^{n-k+1} + \frac{1}{B(k, n-k+1)} \int_0^{F(x)} t^{k-1}(1-t)^{n-k} dt. \\ &= \binom{n}{k-1} F^{k-1}(x)(1-F(x))^{n-k+1} + \sum_{j=k}^n \binom{n}{j} F^j(x)(1-F(x))^{n-j} \\ &= \sum_{j=k-1}^n \binom{n}{j} F^j(x)(1-F(x))^{n-j}. \end{aligned}$$

□

Důsledky.

1. Mají-li X_i rovnoměrné rozdělení na intervalu $(0, 1)$, pak náhodná veličina $X_{(k)}$ má rozdělení s hustotou $f(x) = \frac{1}{B(k, n-k+1)} x^{k-1}(1-x)^{n-k} \mathbb{1}\{x \in (0, 1)\}$, tj. má tzv. beta rozdělení* $B(k, n-k+1)$. Z toho (mimo jiné) plyne

$$\mathbb{E} X_{(k)} = \frac{k}{n+1}, \quad \text{var}(X_{(k)}) = \frac{k(n-k+1)}{(n+2)(n+1)^2}.$$

2. Nechť mají X_i jakékoli spojitě rozdělení s rostoucí distribuční funkcí F . Potom $F(X_{(k)}) \sim B(k, n-k+1)$.

Na druhou stranu nechť $Z \sim B(k, n-k+1)$. Pak

$$\mathbb{P}[X_{(k)} \leq x] = \mathbb{P}[F(X_{(k)}) \leq F(x)] = \mathbb{P}[Z \leq F(x)] = \mathbb{P}[F^{-1}(Z) \leq x],$$

tj. $X_{(k)}$ má stejné rozdělení jako $F^{-1}(Z)$.

Věta 2.14 Hustota k -té pořádkové statistiky vzhledem k Lebesgueově míře jest

$$f_{(k)}(x) = n \binom{n-1}{k-1} f(x) F^{k-1}(x) [1-F(x)]^{n-k}.$$

Důkaz. S využitím věty 2.13

$$f_{(k)}(x) = F'_{(k)}(x) = \frac{1}{B(k, n-k+1)} f(x) F^{k-1}(x) (1-F(x))^{n-k}$$

a tvrzení věty plyne z toho, že

$$\frac{1}{B(k, n-k+1)} \stackrel{(A.1)}{=} \frac{\Gamma(n+1)}{\Gamma(k)\Gamma(n-k+1)} = \frac{n!}{(k-1)!(n-k)!} = \frac{n(n-1)!}{(k-1)!(n-k)!} = n \binom{n-1}{k-1}.$$

□

* Viz např. kapitola 8.2.7 v Kulich (2018).

Věta 2.15 Náhodný vektor $\mathbf{R} = (R_1, \dots, R_n)^\top$ nabývá všech hodnot na množině \mathcal{P}_n , přičemž každá z nich má pravděpodobnost $1/n!$.

Důkaz. Přednášková verze důkazu.

Uvažujme nejdříve pro jednoduchost $n = 2$. Potom

$$P(R_1 = 1, R_2 = 2) = P(X_1 < X_2) = P(X_2 < X_1) = P(R_1 = 2, R_2 = 1),$$

kde jsme podobně jako v důkazu věty 2.12 využili toho, že (díky nezávislosti a stejnému rozdělení X_1 a X_2) má náhodný vektor $(X_1, X_2)^\top$ stejné rozdělení jako náhodný vektor $(X_2, X_1)^\top$. Pro $n = 2$ tedy dostáváme

$$P(R_1 = 1, R_2 = 2) = P(R_1 = 2, R_2 = 1) = \frac{1}{2!}.$$

Výše uvedená úvaha se dá zobecnit pro $n > 2$ následovně. Nechť (i_1, \dots, i_n) je libovolná permutace čísel $(1, \dots, n)$. Potom

$$P(X_1 < X_2 < \dots < X_n) = P(X_{i_1} < X_{i_2} < \dots < X_{i_n})$$

a tudíž

$$P(\mathbf{R}(\mathbf{X}) = \mathbf{r}) = P(\mathbf{R}(\mathbf{X}) = (1, 2, \dots, n)^\top)$$

pro všechny $\mathbf{r} \in \mathcal{P}_n$. Tvrzení věty pak plyne z toho, že množina \mathcal{P}_n má právě $n!$ prvků.

Skriptová verze důkazu.

$$\begin{aligned} P(\mathbf{R}(\mathbf{X}) = \mathbf{r}) &= \int \dots \int \mathbb{1}\{\mathbf{R}(\mathbf{x}) = \mathbf{r}\} f(x_1) \dots f(x_n) dx_1 \dots dx_n \\ &= \int \dots \int \mathbb{1}\{\mathbf{R}(y_{r_1}, \dots, y_{r_n}) = \mathbf{r}\} f(y_{r_1}) \dots f(y_{r_n}) dy_1 \dots dy_n \\ &= \int \dots \int \mathbb{1}\{y_1 < \dots < y_n\} f(y_{r_1}) \dots f(y_{r_n}) dy_1 \dots dy_n \\ &= \int \dots \int \mathbb{1}\{y_1 < \dots < y_n\} f(y_1) \dots f(y_n) dy_1 \dots dy_n \\ &= P(\mathbf{R}(\mathbf{X}) = (1, 2, \dots, n)^\top), \end{aligned}$$

kde jsme podobně jako v důkazu věty 2.12 přeznačili $\mathbf{y} = \mathbf{x}_{(\cdot)}$ (tj. $x_i = y_{r_i}$, $i = 1, \dots, n$), z čehož plyne, že $y_1 < \dots < y_n$.

Z výše uvedeného vyplývá, že

$$P(\mathbf{R} = \mathbf{r}) = \text{const.}, \quad \text{pro } \forall \mathbf{r} \in \mathcal{P}_n.$$

Tvrzení věty pak plyne z toho, že množina \mathcal{P}_n má právě $n!$ prvků. □

Věta 2.16 Platí

- (i) $P(R_i = k) = \frac{1}{n}$ pro všechna $i, k \in \{1, \dots, n\}$.
(ii) $P(R_i = k, R_j = m) = \frac{1}{n(n-1)}$ pro všechna $i \neq j, k \neq m \in \{1, \dots, n\}$.
(iii) $E R_i = \frac{n+1}{2}$, $\text{var } R_i = \frac{n^2-1}{12}$ pro všechna $i \in \{1, \dots, n\}$.
(iv) $\text{cov}(R_i, R_j) = -\frac{n+1}{12}$ pro všechna $i \neq j \in \{1, \dots, n\}$.

Důkaz. Část (i). Bez újmy na obecnosti můžeme uvažovat $i = n$. Dále nechť \mathcal{P}_{n-1}^k obsahuje ty prvky \mathcal{P}_n , které mají na posledním místě číslo k . Nyní

$$P(R_n = k) = \sum_{r \in \mathcal{P}_{n-1}^k} P(\mathbf{R} = r) = (n-1)! \frac{1}{n!} = \frac{1}{n},$$

kde jsme využili větu 2.15 a toho, že množina \mathcal{P}_{n-1}^k má $(n-1)!$ prvků.

Část (ii). Bez újmy na obecnosti můžeme uvažovat $i = n-1$ a $j = n$. Dále nechť $\mathcal{P}_{n-2}^{k,m}$ obsahuje ty prvky \mathcal{P}_n , které mají na předposledním místě číslo k a na posledním místě číslo m . Potom

$$P(R_{n-1} = k, R_n = m) = \sum_{r \in \mathcal{P}_{n-2}^{k,m}} P(\mathbf{R} = r) = (n-2)! \frac{1}{n!} = \frac{1}{n(n-1)},$$

kde jsme využili větu 2.15 a toho, že množina $\mathcal{P}_{n-2}^{k,m}$ má $(n-2)!$ prvků.

Část (iii). Dle části (i):

$$E R_i = \sum_{k=1}^n k P(R_i = k) = \sum_{k=1}^n k \frac{1}{n} = \frac{1}{n} \frac{n(n+1)}{2} = \frac{n+1}{2}.$$

Podobně

$$\begin{aligned} \text{var } R_i &= E R_i^2 - (E R_i)^2 = \sum_{k=1}^n k^2 \frac{1}{n} - \left(\frac{n+1}{2}\right)^2 = \frac{n(n+1)(2n+1)}{6n} - \frac{(n+1)^2}{4} \\ &= \frac{n+1}{12} (4n+2-3n-3) = \frac{(n+1)(n-1)}{12}. \end{aligned}$$

Část (iv).

$$\begin{aligned} \text{cov}(R_i, R_j) &= E R_i R_j - E R_i E R_j = \sum_{k=1}^n \sum_{m=1, m \neq k}^n k m \frac{1}{n(n-1)} - \left(\frac{n+1}{2}\right)^2 \\ &= \frac{1}{n(n-1)} \left[\sum_{k=1}^n k \sum_{m=1}^n m - \sum_{k=1}^n k^2 \right] - \left(\frac{n+1}{2}\right)^2 \\ &= \frac{1}{n(n-1)} \left[\left(\frac{n(n+1)}{2}\right)^2 - \frac{n(n+1)(2n+1)}{6} \right] - \left(\frac{n+1}{2}\right)^2 \end{aligned}$$

$$\begin{aligned}
 &= \frac{n(n+1)^2}{4(n-1)} - \frac{(n+1)(2n+1)}{6(n-1)} - \frac{(n+1)^2}{4} \\
 &= \frac{(n+1)}{12(n-1)} \left[3n(n+1) - 2(2n+1) - 3(n+1)(n-1) \right] \\
 &= \frac{(n+1)}{12(n-1)} (1-n) = -\frac{(n+1)}{12}.
 \end{aligned}$$

□

Poznámka. Pokud data nepocházejí ze spojitého rozdělení nebo se v nich nacházejí shodná pozorování vzniklá vlivem zaokrouhlování, tak dává stále smysl definovat uspořádaný náhodný výběr jako

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n-1)} \leq X_{(n)},$$

přičemž pořádková statistika $X_{(k)}$ je stále dobře definována a platí pro ni věta 2.13.

Pořadí však již nelze stanovit jednoznačně. V takovém případě se často používají tzv. *průměrná pořadí**, která lze spočítat jako

$$\tilde{R}_i = 1 + \sum_{j=1}^n \mathbb{1}\{X_i > X_j\} + \frac{1}{2} \sum_{j=1, j \neq i}^n \mathbb{1}\{X_i = X_j\}. \quad (2.6)$$

Pro takto upravená pořadí však z výše uvedených tvrzení platí pouze $E \tilde{R}_i = \frac{n+1}{2}$ (viz věta 2.16(iii)).

Alternativně je možné pořadí shodných pozorování přiřadit náhodně. Pro náhodně určená pořadí pak platí věta 2.15 a tudíž i věta 2.16. Jelikož se však pořadí většinou používají pro testování, tak výsledek testu by při použití tohoto přístupu mohl záviset na počátečním náhodném přiřazení pořadí ke shodným pozorováním. To se však v praxi považuje za nežádoucí, protože do našich závěrů bychom vnášeli dodatečnou náhodu.

2.4. TRANSFORMACE VE STATISTICE

2.4.1. TRANSFORMACE POZOROVÁNÍ A JEJÍ VLIV NA PARAMETRY

Mějme náhodný výběr X_1, \dots, X_n z rozdělení s distribuční funkcí F_X , hustotou f_X a nosičem S_X . Uvažujme ryze monotonní[†] diferencovatelnou funkci $g : S_X \rightarrow \mathbb{R}$ a definujme $Y_i = g(X_i)$. Potom Y_1, \dots, Y_n je náhodný výběr z rozdělení s hustotou f_Y . Kdyby rozdělení F_X bylo spojitě a kdybychom znali f_X , spočítali bychom hustotu f_Y z tvrzení P.5.3.

Transformace pozorování se ve statistice používají dosti často. Běžný důvod pro provedení transformace bývá, že původní náhodný výběr X_1, \dots, X_n příliš porušuje

* Angl. *average ranks* † Nemonotonním transformacím se obvykle vyhýbáme, protože by mohly ztotožnit pozorování, která byla původně výrazně odlišná.

předpoklady metod, které bychom chtěli použít (například normalitu, symetrii hustoty, existenci momentů apod.). Najdeme tedy vhodnou funkci g takovou, že $Y_i = g(X_i)$ splňuje předpoklady lépe než původní pozorování a pracujeme s náhodným výběrem Y_1, \dots, Y_n namísto původního náhodného výběru X_1, \dots, X_n . Mezi nejčastěji používané transformace kladných náhodných veličin patří např. $g(x) = \log x$ nebo $g(x) = \sqrt{x}$.

Příklad. Nechť X_i má tzv. *logaritmicko-normální rozdělení* $\text{LN}(\mu, \sigma^2)$. Potom $\log(X_i)$ má normální rozdělení $\text{N}(\mu, \sigma^2)$

Pokud používáme transformace, musíme si uvědomovat, že **řada parametrů rozdělení** F_X původního náhodného výběru **se po transformaci změní** takovým způsobem, že je už **nedokážeme identifikovat**.

Například střední hodnota $\mu_X = E X_i$ se změní na $\mu_Y = E g(X_i)$. Pokud neznáme rozdělení X_i , nemůžeme pak z μ_Y spočítat původní střední hodnotu μ_X , ledaže by funkce g byla lineární. Nechť je g rostoucí a ryze konkávní funkce, pak platí z Jensenovy nerovnosti (věta P.2.5) $\mu_Y < g(\mu_X)$ a zpětná transformace $g^{-1}(\mu_Y)$ dává hodnotu ostře menší než μ_X . U ryze konvexní funkce je tomu naopak.

Spočítáme-li tedy výběrový průměr \bar{Y}_n z transformovaného náhodného výběru, bude konvergovat (v pravděpodobnosti) podle věty 2.2(ii) k μ_Y . Zpětná transformace $g^{-1}(\bar{Y}_n)$ bude konvergovat (v pravděpodobnosti) k $g^{-1}(\mu_Y) \neq \mu_X$. Obecně nelze nalézt funkci h takovou, aby $h(\bar{Y}_n)$ konvergovalo k μ_X . Zajímá-li nás konkrétní hodnota μ_X , nemůžeme tedy data transformovat. Podobné je to s rozptylem a vyššími momenty: po transformaci už obvykle nezjistíme, jaký byl rozptyl původních pozorování.

Příklad. Nechť $X_i \sim \text{LN}(\mu, \sigma^2)$. Potom pro $g(x) = \log x$ platí, že $Y_i = g(X_i) \sim \text{N}(\mu, \sigma^2)$. Tedy

$$g^{-1}(\bar{Y}_n) \xrightarrow[n \rightarrow \infty]{P} e^{E Y_i} = e^\mu < e^{\mu + \sigma^2/2} = E X_i.$$

Některé jiné parametry však tento problém nemají. Například medián nebo kterýkoli jiný kvantil lze snadno získat zpětnou transformací: Nechť m_X je medián X_i a m_Y je medián Y_i , nechť g je rostoucí funkce. Pak platí $m_Y = g(m_X)$, tj. m_X lze identifikovat zpětnou transformací $g^{-1}(m_Y)$.

Pořadí jsou invariantní vůči rostoucím transformacím, takže statistiky závislé pouze na pořadích nabývají stejné hodnoty, ať už jsou počítány z původního nebo transformovaného náhodného výběru.

2.4.2. TRANSFORMACE STABILIZUJÍCÍ (ASYMPTOTICKÝ) ROZPTYL

Jinou motivací pro použití transformace může být snaha stabilizovat (asymptotický) rozptyl. Mějme posloupnost náhodných veličin T_n , které splňují, že

$$\sqrt{n} (T_n - \mu) \xrightarrow[n \rightarrow \infty]{d} \text{N}(0, \sigma^2(\mu)).$$

Rozptyl $\sigma^2(\mu)$ asymptotického normálního rozdělení se někdy nazývá také asymptotický rozptyl* posloupnosti $\sqrt{n}(T_n - \mu)$.

Jak uvidíme později, pro inferenci (testování, intervaly spolehlivosti) o parametru μ je zpravidla dobré, pokud asymptotický rozptyl již nezávisí na parametru μ .

Nechť tedy g je nějaká reálná funkce, která je definovaná a diferencovatelná na okolí bodu μ . Potom pomocí Δ -metody (Tvzení 1.7) dostáváme, že

$$\sqrt{n}(g(T_n) - g(\mu)) \xrightarrow[n \rightarrow \infty]{d} N(0, [g'(\mu)]^2 \sigma^2(\mu)).$$

Pokud tedy budeme volit

$$g(x) = c \int \frac{1}{\sigma(x)} dx, \quad (2.7)$$

potom $g'(\mu) = \frac{c}{\sigma(\mu)}$ a tudíž

$$\sqrt{n}(g(T_n) - g(\mu)) \xrightarrow[n \rightarrow \infty]{d} N(0, c^2)$$

a vliv μ na asymptotický rozptyl bude eliminován.

Příklad. Nechť X_1, \dots, X_n je náhodný výběr z Poissonova rozdělení $Po(\lambda)$. Potom statistika $T_n = \bar{X}_n$ dle centrální limitní věty (tvrzení 1.6) splňuje

$$\sqrt{n}(\bar{X}_n - \lambda) \xrightarrow[n \rightarrow \infty]{d} N(0, \lambda).$$

Tedy $\sigma(x) = \sqrt{x}$ a tudíž $g(x) = \int \frac{1}{\sigma(x)} dx = \int x^{-1/2} dx = 2\sqrt{x}$ a dostáváme, že

$$\sqrt{n}(2\sqrt{\bar{X}_n} - 2\sqrt{\lambda}) \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Poznámka. Podobná myšlenka se někdy využívá i pro samostatná pozorování. Nechť platí $E X_i = \lambda$ a $\text{var } X_i = \sigma^2(\lambda)$. Potom doufáme, že po přechodu k transformaci $Y_i = g(X_i)$, kde g se spočte pomocí (2.7), budou mít pozorování Y_i rozdělení bližší normálnímu. Tedy např. pro $X_i \sim Po(\lambda)$ se často pracuje s $Y_i = \sqrt{X_i}$.

Cvičení. Uvažujte posloupnost náhodných veličin $\{Y_n\}$ takovou že, $Y_n \sim Po(n\lambda)$. Dokažte, že $\sqrt{Y_n} - \sqrt{n\lambda} \xrightarrow[n \rightarrow \infty]{d} N(0, \frac{1}{4})$.

Nápověda: Uvědomte si, že Y_n se dá reprezentovat jako $\sum_{i=1}^n X_i$, kde X_1, \dots, X_n je náhodný výběr z $Po(\lambda)$.

2.4.3. STANDARDIZACE

Speciálním druhem transformace je tzv. *standardizace*. Máme náhodný výběr X_1, \dots, X_n a spočítáme \bar{X}_n a S_n^2 . Potom definujeme náhodné veličiny Z_1, \dots, Z_n vztahem

$$Z_i = \frac{X_i - \bar{X}_n}{S_n}.$$

* Angl. *asymptotic variance*

Tyto veličiny mají výběrový průměr 0 a výběrový rozptyl 1, ale nepředstavují náhodný výběr, neboť nejsou nezávislé. Jelikož však

$$\bar{X}_n \xrightarrow[n \rightarrow \infty]{P} E X_i, \quad \text{a také} \quad S_n \xrightarrow[n \rightarrow \infty]{P} \sqrt{\text{var } X_i},$$

tak při dostatečně velkém počtu pozorování se Z_1, \dots, Z_n chovají téměř jako nezávislé veličiny s nulovou střední hodnotou a jednotkovým rozptylem. V mnoha případech lze ukázat, že závislost vznikla tím, že jsme neznámé $E X_i$ a $\sqrt{\text{var } X_i}$ nahradili jejich výběrovými protějšky (tj. \bar{X}_n a S_n) lze zanedbat.

Standardizace se používá tehdy, pokud se chceme zbavit prvních dvou momentů a soustředit se na jiné aspekty rozdělení F_X (viz např. výběrový korelační koeficient v Kapitole 10.1).

Přípravné příklady ke zkoušce.

1. Nechť X_1 a X_2 jsou nezávislé náhodné veličiny s rovnoměrným rozdělením na intervalu $(0, 1)$. Spočtete $E \frac{X_1^2}{X_2}$ a $E \frac{X_1^2}{X_{(2)}}$.
2. Nechť X_1, X_2, X_3 jsou nezávislé náhodné veličiny s rovnoměrným rozdělením na intervalu $(0, 1)$. Spočtete $E \frac{X_3}{X_2}$ a $E \frac{X_{(3)}}{X_{(2)}}$.
3. Nechť X_1, \dots, X_n jsou nezávislé stejně rozdělené náhodné veličiny s hustotou f vzhledem k Lebesgueově míře. Spočtete pravděpodobnost $P(R_1 + R_2 = 6)$.
4. Nechť X_1, \dots, X_n jsou nezávislé stejně rozdělené náhodné veličiny s hustotou f vzhledem k Lebesgueově míře. Spočtete $E \frac{R_1}{R_2}$ a následně $\lim_{n \rightarrow \infty} E \frac{R_1}{R_2}$.

3. ODHADOVÁNÍ PARAMETRŮ

Máme náhodný výběr $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$, model \mathcal{F} a parametr $\theta = t(F) \in \mathbb{R}^p$ pro $F \in \mathcal{F}$, který chceme v daném modelu odhadnout. Nechť $F_X \in \mathcal{F}$ je skutečné rozdělení náhodného vektoru \mathbf{X}_i a $\theta_X \equiv t(F_X)$ je skutečná hodnota hledaného parametru.

3.1. BODOVÝ ODHAD

Definice 3.1 *Odhadem parametru $\theta_X \equiv t(F_X) \in \mathbb{R}^p$ rozumíme p -rozměrný náhodný vektor $\hat{\theta}_n$, který spočteme jako $\hat{\theta}_n = T_n(\mathbf{X}) \equiv T_n(\mathbf{X}_1, \dots, \mathbf{X}_n)$, kde T_n je nějaká borelovsky měřitelná funkce dat.**

Poznámka. Odhad je statistika ve smyslu definice 2.3. Odhad nesmí záviset na neznámých parametrech.

Definice 3.2 (Nestrannost a konsistence) Mějme náhodný výběr $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$ z rozdělení $F_X \in \mathcal{F}$ a odhad $\hat{\theta}_n \equiv T_n(\mathbf{X})$ parametru $\theta_X \equiv t(F_X)$.

- (i) Řekneme, že odhad $\hat{\theta}_n$ je *nestranný odhad*[†] parametru θ_X v modelu \mathcal{F} , právě když $E \hat{\theta}_n = \theta_X$ pro každé n (pro něž je odhad definován) a pro každé rozdělení $F_X \in \mathcal{F}$.
- (ii) Řekneme, že odhad $\hat{\theta}_n$ je *konsistentní odhad*[‡] parametru θ_X v modelu \mathcal{F} , právě když $\hat{\theta}_n \xrightarrow{P} \theta_X$ při $n \rightarrow \infty$ pro každé rozdělení $F_X \in \mathcal{F}$.

Poznámka.

- Vlastnosti odhadů musíme zkoumat v kontextu daného modelu. Snadno se může stát, že odhad $\hat{\theta}_n$ je nestranný a konsistentní v nějakém modelu \mathcal{F} , ale v jiném modelu \mathcal{F}' tyto vlastnosti nemá.
- Nestrannost má platit pro každý počet pozorování n , pro něž je odhad definován (např. u výběrového rozptylu pro $n \geq 2$). Nestrannost ale nezaručuje, že se odhad při zvětšujícím se rozsahu výběru přibližuje k hledanému parametru. Pro některé modely neexistují rozumné (nebo vůbec žádné) nestranné odhady.
- Konsistence je asymptotická vlastnost, která nic neříká o chování odhadu při konečném n . (Příklad: $\hat{\theta}_n = 21,5$ pro $n \leq 10^{10}$, $\hat{\theta}_n = \bar{X}_n$ pro $n > 10^{10}$ je konsistentní odhad $\theta_X = E X_i$.)
- Námi definovaná konzistence se někdy také nazývá *slabá konzistence*[§]. Odhad se pak nazývá *silně konsistentní*[¶], pokud platí $\hat{\theta}_n \xrightarrow[n \rightarrow \infty]{s.j.} \theta_X$.

* Angl. *estimator, estimate* † Angl. *unbiased estimator* ‡ Angl. *consistent estimator* § Angl. *weak consistency* ¶ Angl. *strong consistency*

- Odhady, které nejsou nestranné, ale jsou konsistentní, se ve statistice běžně používají. Odhady, které nejsou konsistentní, nepoužíváme, neboť odhadují „něco jiného“ nebo se s rostoucím rozsahem výběru „nezpřesňují“.

Příklady.

1. *Odhad parametru* $\theta_X = E X_i$ *v modelu* $\mathcal{F} = \mathcal{L}^1$:
 - Průměr \bar{X}_n je nestranný a konsistentní odhad θ_X [plyne z věty 2.2, (i) a (ii)].
 - Odhad $\hat{\theta}_n = X_1$ je nestranný odhad θ_X , ale není konsistentní.
2. *Odhad parametru* $\theta_X = \text{var } X_i$ *v modelu* $\mathcal{F} = \mathcal{L}^2$:
 - Výběrový rozptyl S_n^2 je nestranný a konsistentní odhad θ_X [plyne z věty 2.6, (i) a (ii)].
 - Odhad $\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ je konsistentní odhad θ_X , ale není nestranný.
3. *Odhad parametru* $\theta_X = P[X_i = 0]$ *v modelu* $\mathcal{F} = \{\text{Po}(\lambda), \lambda > 0\}$:
 - Odhad $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i = 0\}$ je nestranný a konsistentní odhad θ_X (a to dokonce v modelu všech diskretních rozdělení).
 - Odhad $\tilde{\theta}_n = \left(\frac{n-1}{n}\right)^{\sum_{i=1}^n X_i}$ je také nestranný a konsistentní odhad θ_X (v modelu \mathcal{F} nikoliv však v modelu všech diskretních rozdělení).
4. *Odhad parametru* $\theta_X = e^{-2\lambda x}$ *v modelu* $\mathcal{F} = \{\text{Po}(\lambda), \lambda > 0\}$ *pro* $n = 1$:
 Jediný nestranný odhad jest $\hat{\theta} = (-1)^{X_1}$, jeho možné hodnoty jsou -1 a 1 . Hledaný parametr $e^{-2\lambda x}$ však nabývá pouze hodnot z intervalu $(0, 1)$.

Definice 3.3 (Vychýlení) Nechť odhad $\hat{\theta}_n \equiv T_n(\mathbf{X})$ parametru θ_X má konečnou střední hodnotu. Rozdíl $E(\hat{\theta}_n - \theta_X)$ nazýváme *vychýlením** odhadu $\hat{\theta}_n$.

Definice 3.4 Nechť odhad $\hat{\theta}_n \equiv T_n(\mathbf{X})$ parametru $\theta_X \in \mathbb{R}$ má konečný rozptyl.

(i) Výraz

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

nazýváme *střední čtvercovou chybou* odhadu $\hat{\theta}_n$.[†]

(ii) Výraz

$$\text{SE}(\hat{\theta}_n) = \sqrt{\text{var}(\hat{\theta}_n)}$$

nazýváme *směrodatnou chybou*[‡] odhadu $\hat{\theta}_n$.

Poznámka.

* Angl. *bias* † Angl. *mean square error, MSE* ‡ Angl. *standard error, SE*

3. Odhadování parametrů

- Pozor na jemné rozdíly v terminologii. Pojem *směrodatná odchylka* (standard deviation, SD) obvykle znamená odmocninu z rozptylu jednoho pozorování náhodného výběru, tj. $\sqrt{\text{var } X_i}$. Pojem *směrodatná chyba* (standard error, SE) obvykle znamená odmocninu z rozptylu nějakého odhadu spočítaného z celého náhodného výběru. Někteří autoři však pojmem *směrodatná chyba* rozumí

$$\text{SE}(\hat{\theta}_n) = \sqrt{\widehat{\text{var}}(\hat{\theta}_n)},$$

kde $\widehat{\text{var}}(\hat{\theta}_n)$ je odhad $\text{var}(\hat{\theta}_n)$

- Střední čtvercová chyba i směrodatná chyba jsou míry *přesnosti* odhadu. Směrodatná chyba do přesnosti nezahrnuje vychýlení, zatímco střední čtvercová chyba ano.
- Platí, že střední čtvercová chyba lze rozložit na rozptyl a kvadrát vychýlení, tj. :

$$\text{MSE}(\hat{\theta}_n) = \text{var}(\hat{\theta}_n) + [\text{E}(\hat{\theta}_n - \theta_X)]^2.$$

Důkaz výše uvedeného rozkladu plyne z toho, že

$$\begin{aligned} \text{MSE}(\hat{\theta}_n) &= \text{E}(\hat{\theta}_n - \text{E}\hat{\theta}_n + \text{E}\hat{\theta}_n - \theta_X)^2 \\ &= \text{E}(\hat{\theta}_n - \text{E}\hat{\theta}_n)^2 + 2\text{E}(\hat{\theta}_n - \text{E}\hat{\theta}_n)\text{E}(\hat{\theta}_n - \theta_X) + [\text{E}(\hat{\theta}_n - \theta_X)]^2 \\ &= \text{var}(\hat{\theta}_n) + 0 + [\text{E}(\hat{\theta}_n - \theta_X)]^2. \end{aligned}$$

- Střední čtvercová chyba je jedno z nevhodnějších kritérií pro porovnávání odhadů. Máme-li několik různých odhadů téhož parametru v tomtéž modelu, snažíme se mezi nimi najít ten, který má nejmenší MSE. Tj. v případě nestranných odhadů vybíráme odhad s nejmenším rozptylem.
- MSE často nelze spočítat. V mnoha případech se však lze rozhodovat na základě asymptotického rozptylu odhadů. Tj. předpokládejme, že máme dva odhady $\hat{\theta}_n$ a $\tilde{\theta}_n$, které splňují

$$\sqrt{n}(\hat{\theta}_n - \theta_X) \xrightarrow[n \rightarrow \infty]{d} \text{N}(0, \sigma_1^2), \quad \sqrt{n}(\tilde{\theta}_n - \theta_X) \xrightarrow[n \rightarrow \infty]{d} \text{N}(0, \sigma_2^2).$$

Potom (pro velké rozsahy výběrů) preferujeme odhad $\hat{\theta}_n$ pokud $\sigma_1^2 < \sigma_2^2$ nebo naopak odhad $\tilde{\theta}_n$ pokud $\sigma_1^2 > \sigma_2^2$.

Příklad. Odhad parametru $\sigma_X^2 = \text{var } X_i$ v modelu $\mathcal{F} = \{\text{N}(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$. Ukažte, že platí: $\text{MSE}(S_n^2) > \text{MSE}(\hat{\sigma}_n^2)$.

Věta 3.1 Nechť $\hat{\theta}_n$ je odhad parametru $\theta_X \in \mathbb{R}$, pro nějž platí

$$\text{E}\hat{\theta}_n \xrightarrow[n \rightarrow \infty]{} \theta_X \quad (\text{vychýlení konverguje k nule}) \quad \text{a také} \quad \text{var}(\hat{\theta}_n) \xrightarrow[n \rightarrow \infty]{} 0$$

pro všechna $F_X \in \mathcal{F}$. Pak je $\hat{\theta}_n$ konsistentní odhad θ_X .

Důkaz. Necht' $\varepsilon > 0$. Potom s využitím předpokladů věty a Markovovy nerovnosti (věty P.2.6):

$$P(|\hat{\theta}_n - \theta_X| > \varepsilon) \leq \frac{\text{MSE}(\hat{\theta}_n)}{\varepsilon^2} = \frac{\text{var}(\hat{\theta}_n)}{\varepsilon^2} + \frac{(E\hat{\theta}_n - \theta_X)^2}{\varepsilon^2}.$$

Nyní první i druhý člen na pravé straně k nule, protože dle předpokladu $\text{var}(\hat{\theta}_n) \rightarrow 0$ a $E\hat{\theta}_n \rightarrow \theta_X$ pro $n \rightarrow \infty$. \square

Poznámka.

- Opačná implikace neplatí. Existují běžně používané konsistentní odhady, pro něž platí $E|\hat{\theta}_n| = \infty$ pro každé konečné n .
- Věta 3.1 je šikovná v situacích, kdy máme k dispozici (či lze snadno spočítat) vychýlení a rozptyl odhadu $\hat{\theta}_n$. Pokud však můžeme psát $\hat{\theta}_n = g(\frac{1}{n} \sum_{i=1}^n X_i)$ (tj. jako transformaci výběrového průměru), pak lze konzistence vyšetřovat jednodušeji kombinací zákona velkých čísel (tvrzení 1.5) a věty o spojitě transformaci (tvrzení 1.2).

Příklad. Necht' X_1, \dots, X_n je náhodný výběr z alternativního rozdělení $\text{Alt}(p_X)$. Uvažujte $\hat{\theta}_n = \frac{1}{X_n}$ jako odhad parametru $\theta_X = \frac{1}{p_X}$. Ukažte, že přestože $E\hat{\theta}_n = \infty$, tak $\hat{\theta}_n \xrightarrow[n \rightarrow \infty]{P} \theta_X$.

3.2. VOLBA PARAMETRU

Parametr $\theta = t(F)$, který se snažíme odhadovat, může být v principu cokoli. Ne všechny parametry však dávají smysl v kontextu daného praktického problému, který řešíme. Musíme tedy rozlišovat, které parametry pro daný problém má smysl odhadovat a které ne. To záleží na významu hodnot měřených veličin, na tom, jak byly získány, zpracovány atd. Statistické metody, kterými se budeme zabývat, budeme rozlišovat podle toho, pro jaký typ měření jsou určeny. Přitom budeme uvažovat následující typy dat, neboli *škály měření**.

3.2.1. KVANTITATIVNÍ DATA

Náhodnou veličinu X nazveme *kvantitativní†*, pokud její hodnoty mají konkrétní numerický význam (např. počet, procento, délka, objem, hmotnost, úroková míra, koncentrace látky, energie, teplota, doba trvání, velikost úhlu, zeměpisná šířka, kalendářní rok). U kvantitativních veličin existuje smysluplné uspořádání jejich hodnot (teplota 10 °C je vyšší než -11,4 °C) a rozdíly jejich hodnot mají reálnou interpretaci. Kvantitativní veličiny mohou být jak diskrétní tak spojitě.

Kvantitativní veličiny můžeme dále dělit na dvě podskupiny: *intervalové* a *poměrové*. **Poměrové veličiny** jsou typicky nezáporné s jasně definovanou nulovou hodnotou a interpretovatelnými podíly. Například hmotnost 0 kg je jednoznačně daná

* Angl. *measurement scales* † Angl. *quantitative*

a hmotnost 20 kg je čtyřikrát více než 5 kg. Příklady poměrových veličin jsou počet, délka, objem, hmotnost, úroková míra, koncentrace látky, energie, doba trvání, teplota měřená v Kelvinech. **Intervalové veličiny** jsou kvantitativní veličiny, které nejsou poměrové, to jest nemají pevně definovanou nulu nebo nemají interpretovatelné podíly. Například směr daný azimutem je intervalová veličina, neboť azimut 360° není šestkrát větší než 60° . Podobně teplota měřená v $^\circ\text{C}$ je intervalová veličina neboť 16°C není čtyřikrát vyšší teplota než 4°C . Kalendářní rok je také intervalová veličina, protože nemá smysl počítat podíl letošního roku a roku vašeho narození.

3.2.2. KATEGORIÁLNÍ DATA

Náhodnou veličinu X nazveme *kategoriální**, pokud její hodnoty kódují příslušnost (neboli *klasifikaci*) subjektu do určité kategorie, neboli jedné z několika disjunktních množin. Kategoriální veličiny jsou vždy diskrétní a mají konečný počet K možných hodnot, obvykle $1, \dots, K$ nebo $0, \dots, K - 1$. Hodnoty kategoriálních veličin nemají přímou numerickou interpretaci, slouží pouze k rozlišení konečného počtu možných stavů. Jednotlivým stavům říkáme *úrovně†* nebo *kategorie*.

Kategoriální veličiny dále dělíme na *nominální‡* a *ordinální§*. U **nominálních veličin** neexistuje ani žádné uspořádání jejich kategorií – nelze říci, že kategorie j předchází kategorii $j + 1$. Příkladem nominální veličiny je třeba bydliště kategorizované jako kraj (1 = Praha, 2 = Středočeský kraj, ..., 14 = Zlínský kraj) nebo sociální postavení (1 = nezletilý; 2 = student; 3 = zaměstnanec; 4 = živnostník; 5 = nezaměstnaný; 6 = důchodce). **Ordinální veličiny** mají v nějakém smyslu uspořádané kategorie, takže lze tvrdit, že kategorie j předchází kategorii $j + 1$, nebo že je menší, horší apod. Příkladem ordinální veličiny je třeba odpověď na otázku s možnostmi 1 = ostře nesouhlasím, 2 = spíše nesouhlasím, 3 = nevím, 4 = spíše souhlasím, 5 = naprosto souhlasím. Jiný příklad je veličina nejvyšší dosažené vzdělání kódovaná jako 1 = nižší než základní; 2 = základní; 3 = učební obor; 4 = středoškolské s maturitou; 5 = bakalářské; 6 = magisterské; 7 = doktorské.

3.2.3. BINÁRNÍ DATA

Binární¶ veličiny jsou speciálním případem kategoriálních veličin, kde $K = 2$. Klasifikují tedy pozorování do jednoho ze dvou možných stavů. Jejich hodnoty se obvykle volí jako 0 vs. 1, případně 1 vs. 2. Příkladem binární veličiny je pravdivostní hodnota výroku (0 = pravda, 1 = lež), realizace náhodného jevu (0 = nenastal/neúspěch, 1 = nastal/úspěch) nebo pohlaví (1 = samec, 2 = samice).

3.2.4. VOLBA PARAMETRU V ZÁVISLOSTI NA TYPU DAT

Pro nominální veličiny obecně nemá smysl uvažovat parametry jako $E X$, $\text{var } X$, distribuční funkci, kvantily, kovariance a korelace, zkrátka žádné charakteristiky, které

* Angl. *categorical* † Angl. *levels* ‡ Angl. *nominal* § Angl. *ordinal* ¶ Angl. *binary*

závisí na kódování a uspořádání jednotlivých kategorií. Tyto parametry jsou sice řádně definovány, ale nemají žádnou praktickou interpretaci. Jediné parametry, které u nominálních veličin interpretaci mají, jsou pravděpodobnosti jednotlivých kategorií, čili $p_j = P[X = j]$ pro všechny možné hodnoty j .

Výjimkou jsou binární veličiny. Znamená-li např. hodnota 0 neúspěch a hodnota 1 úspěch, pak $E X = P[X = 1]$, tedy střední hodnota je zároveň pravděpodobnost úspěchu.

U ordinálních veličin má díky uspořádání jejich hodnot smysl distribuční funkce. Často je možné přikládat jim intervalovou interpretaci (doktorské vzdělání je o dva stupně vyšší než bakalářské), ale obvykle jim nelze dávat poměrovou interpretaci (nelze říci, že magisterské vzdělání je dvakrát vyšší než učební obor). Ordinálním veličinám se někdy přiřazují neceločíselné hodnoty, tzv. *skóry*. Např. ordinální veličinu můžeme vytvořit tak, že vezmeme kvantitativní veličinu Z a seskupíme ji podle zvolených dělicích bodů, např. $X = 1$ pokud $Z \in \langle 0, 5 \rangle$, $X = 2$ pokud $Z \in \langle 5, 20 \rangle$, $X = 3$ pokud $Z \in \langle 20, 100 \rangle$ a $X = 4$ pokud $Z \geq 100$. Takové veličiny běžně vznikají v dotaznících, kde respondent dostane na výběr jednu ze čtyř možností namísto toho, aby musel zapsat přesné číslo. Výsledná veličina X je zjevně ordinální. Namísto hodnot 1, ..., 4 bychom ale mohli za hodnoty X vzít prostředky intervalů, z kterých hodnoty X vznikly, tedy 2,5; 12,5 a 60 pro první tři intervaly. S posledním je zjevně potíž, neboť nemá pravý okraj – jeho skóru bychom museli nějak doplnit, například vzít 150. Takto zakódovaná veličina X je nejen ordinální, ale má některé vlastnosti veličiny kvantitativní.

Ordinální veličiny můžeme vždy analyzovat jako by byly nominální, ale často je možné na ně používat metody určené pro kvantitativní veličiny, odhadovat jejich střední hodnotu nebo počítat jejich rozdíly. Existují také speciální metody určené právě pro ordinální veličiny, s těmi se ale zatím nesetkáme.

Náš výklad statistických metod počínaje kapitolou 4 bude rozlišovat metody pro kvantitativní data, kde budeme pracovat s charakteristikami jako je střední hodnota, rozptyl, medián, distribuční funkce, kovariance apod., a metody pro nominální data, kde budeme pracovat s pravděpodobnostmi jednotlivých kategorií.

3.3. MOMENTOVÁ METODA

Momentová metoda* patří spolu s metodou maximální věrohodnosti k základním metodám odhadu parametrů.

Uvažujme nyní parametrický model: máme náhodný výběr X_1, \dots, X_n z rozdělení s hustotou $f(x; \theta_X)$ vůči nějaké σ -konečné míře μ , kde tvar funkce $f(\cdot; \cdot)$ je známý a θ_X je neznámý (vektorový) parametr, jenž leží v parametrickém prostoru $\Theta \subseteq \mathbb{R}^p$, $p \geq 1$. Pracujeme tedy s modelem

$$\mathcal{F} = \{ \text{rozdělení s hustotou } f(x; \theta), \theta \in \Theta \subseteq \mathbb{R}^p \}$$

* Angl. *method of moments*

3. Odhadování parametrů

Cílem je odhadnout parametr θ_X . Využijeme toho, že máme k dispozici konsistentní odhady momentů a že momenty rozdělení X_i obvykle umíme vyjádřit jako funkce neznámých parametrů. Budeme předpokládat, že $E |X_i|^p < \infty$.

Uvažujme nejprve $p = 1$. Předpokládejme, že $E X_i = \tau(\theta_X)$, kde $\tau : \Theta \rightarrow \mathbb{R}$. Jelikož \bar{X}_n je konsistentní odhad, tak se nabízí hledat *momentový odhad** $\hat{\theta}_n$ jako řešení *odhadovací rovnice*†:

$$\bar{X}_n = \tau(\hat{\theta}_n). \quad (3.1)$$

Pokud je funkce τ ryze monotonní, můžeme odhad vyjádřit jako $\hat{\theta}_n = \tau^{-1}(\bar{X}_n)$ a odhadovaný parametr jako $\theta_X = \tau^{-1}(E X_i)$.

Vlastnosti odhadu $\hat{\theta}_n$:

- Je-li τ^{-1} spojitá funkce v bodě $E X_i$, pak $\hat{\theta}_n \xrightarrow[n \rightarrow \infty]{P} \theta_X$ (viz tvrzení 1.2).
- Má-li τ^{-1} spojitou derivaci na okolí bodu $E X_i$, pak pomocí Δ -metody (tvrzení 1.7)

$$\sqrt{n} (\hat{\theta}_n - \theta_X) \xrightarrow[n \rightarrow \infty]{d} N(0, V(\theta_X)),$$

kde

$$V(\theta_X) = \left\{ [\tau^{-1}(E X_i)]' \right\}^2 \text{var } X_i = \frac{\text{var } X_i}{[\tau'(\tau^{-1}(E X_i))]^2} = \frac{\text{var } X_i}{[\tau'(\theta_X)]^2}. \quad (3.2)$$

Povšimněme si, že ve vyjádření asymptotického rozptylu pomocí poslední rovnosti nepotřebujeme znát explicitní předpis pro τ^{-1} . Toto vyjádření se tedy hodí, pokud τ^{-1} je dána pouze implicitně a odhad $\hat{\theta}_n$ hledáme pomocí numerických metod jako řešení odhadovací rovnice (3.1).

V aplikacích asymptotický rozptyl $V(\theta_X)$ odhadujeme pomocí

$$\hat{V}_n = \left\{ [\tau^{-1}(\bar{X}_n)]' \right\}^2 S_n^2 = \frac{S_n^2}{[\tau'(\hat{\theta}_n)]^2},$$

přičemž druhé vyjádření se opět hodí zejména v případě, kdy nemáme explicitní vyjádření pro τ^{-1} .

Příklady.

1. X_1, \dots, X_n je náhodný výběr z rozdělení $\text{Po}(\lambda_X)$, $E X_i = \lambda_X$. Momentovým odhadem parametru λ_X je $\hat{\theta}_n = \bar{X}_n$.
2. X_1, \dots, X_n je náhodný výběr z rozdělení $\text{Geo}(p_X)$, $E X_i = \frac{1-p_X}{p_X}$ a $\text{var } X_i = \frac{1-p_X}{p_X^2}$. Tedy $\tau(x) = \frac{1-x}{x}$ a $\tau^{-1}(x) = \frac{1}{1+x}$. Momentovým odhadem parametru p_X je $\hat{p}_n = \frac{1}{1+\bar{X}_n}$.
Dále

$$\sqrt{n} (\hat{p}_n - p_X) \xrightarrow[n \rightarrow \infty]{d} N(0, p_X^2(1 - p_X)),$$

* Angl. *moment estimator* † Angl. *estimating equation*

3. Odhadování parametrů

kde asymptotický rozptyl $p_X^2(1 - p_X)$ plyne buď z první rovnosti v (3.2)

$$V(p_X) = \left\{ \frac{-1}{(1 + E X_i)^2} \right\}^2 \text{var } X_i = p_X^4 \frac{1 - p_X}{p_X^2}$$

nebo alternativně také z třetí rovnosti v (3.2)

$$V(p_X) = \frac{\text{var } X_i}{\left\{ -\frac{1}{p_X} \right\}^2} = \frac{\frac{1 - p_X}{p_X^2}}{\frac{1}{p_X^4}}.$$

3. X_1, \dots, X_n je náhodný výběr z rozdělení $R(0, \theta_X)$, $E X_i = \theta_X/2$. Momentovým odhadem parametru θ_X je $\hat{\theta}_n = 2\bar{X}_n$. Platí $\sqrt{n}(\hat{\theta}_n - \theta_X) \xrightarrow[n \rightarrow \infty]{d} N(0, \theta_X^2/3)$.

$p = 1$, ale jiný moment než $E X_i$

Někdy se může stát, že $E X_i = 0$ pro všechny $\theta_X \in \Theta$. To platí například pro rozdělení s konečnou střední hodnotou, která jsou symetrická kolem nuly. Potom můžeme uvažovat druhý moment, tj. $E X_i^2 = \tau(\theta_X)$ a odhad $\hat{\theta}_n$ dostaneme jako řešení rovnice

$$\frac{1}{n} \sum_{i=1}^n X_i^2 = \tau(\hat{\theta}_n).$$

Obecně můžeme vzít nějakou šikovnou (měřitelnou) funkci t takovou, že $E |t(X_i)| < \infty$ a $E t(X_i) = \tau(\theta_X)$. Odhad $\hat{\theta}_n$ dostaneme jako řešení rovnice

$$\frac{1}{n} \sum_{i=1}^n t(X_i) = \tau(\hat{\theta}_n).$$

Nyní rozšíříme metodu na $p > 1$.

Nejpřímochařejší je v tomto případě uvažovat prvních p -momentů, tj. spočteme

$$E X_i = \tau_1(\theta_X), E X_i^2 = \tau_2(\theta_X), \dots, E X_i^p = \tau_p(\theta_X),$$

a získáme tak zobrazení $\tau_1, \dots, \tau_p : \Theta \rightarrow \mathbb{R}$. Odhad parametru $\hat{\theta}_n$ pak dostaneme jako řešení soustavy p rovnic o p neznámých

$$\frac{1}{n} \sum_{i=1}^n X_i = \tau_1(\hat{\theta}_n), \frac{1}{n} \sum_{i=1}^n X_i^2 = \tau_2(\hat{\theta}_n), \dots, \frac{1}{n} \sum_{i=1}^n X_i^p = \tau_p(\hat{\theta}_n).$$

Zavedeme-li zobrazení $\tau = (\tau_1, \dots, \tau_p)^\top : \Theta \rightarrow \mathbb{R}^p$, pak za předpokladu existence inverze zobrazení τ můžeme psát

$$\hat{\theta}_n = \tau^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{Z}_i \right), \quad \text{kde } \mathbf{Z}_i = (X_i, X_i^2, \dots, X_i^p)^\top.$$

3. Odhadování parametrů

Z tohoto vyjádření pak podobně jako v případě $p = 1$ můžeme odvodit konzistenci a asymptotickou normalitu odhadu $\widehat{\theta}_n$.

Speciální případ $p = 2$

Předpokládejme, že $(E X_i, \text{var } X_i)^\top = \tau(\theta_X)$, kde $\tau : \Theta \rightarrow \mathbb{R}^2$. Pak se nabízí hledat odhad parametru θ_X jako řešení soustavy odhadovacích rovnic (přesněji dvou rovnic o dvou neznámých)

$$(\overline{X}_n, S_n^2)^\top = \tau(\widehat{\theta}_X).$$

Pokud je funkce τ prostá, tak můžeme odhad vyjádřit jako $\widehat{\theta}_X = \tau^{-1}(\overline{X}_n, S_n^2)$ a odhadovaný parametr jako $\theta_X = \tau^{-1}(E X_i, \text{var } X_i)$.

Vlastnosti odhadu $\widehat{\theta}_n$:

- Víme, že \overline{X}_n a S_n^2 jsou konsistentní odhady $E X_i$ a $\text{var } X_i$. Je-li tedy funkce τ^{-1} spojitá v bodě $(E X_i, \text{var } X_i)$, pak $\widehat{\theta}_n \xrightarrow[n \rightarrow \infty]{P} \theta_X$.
- Z věty 2.6, část (iv) víme, že pokud $E X_i^4 < \infty$, pak \overline{X}_n a S_n^2 jsou sdruženě asymptoticky normální. Má-li τ^{-1} spojitou derivaci, pak podle Δ -metody má i $\widehat{\theta}_n$ asymptoticky sdružené normální rozdělení s rozptylovou maticí, kterou lze spočítat pomocí věty 2.6 a Δ -metody.

Příklady.

4. X_1, \dots, X_n je náhodný výběr z gama rozdělení s parametry s hustotou $f(x; a, p) = \frac{a^p}{\Gamma(p)} x^{p-1} e^{-ax} \mathbb{1}\{x > 0\}$. Pak $E X_i = \frac{p}{a}$ a $\text{var } X_i = \frac{p}{a^2}$ (viz např. kapitola 8.2.6 v Kulich, 2018). Momentovou metodou dostaneme konsistentní a asymptoticky normální odhady

$$\widehat{a}_n = \frac{\overline{X}_n}{S_n^2} \quad \text{a} \quad \widehat{p}_n = \frac{\overline{X}_n^2}{S_n^2}.$$

5. X_1, \dots, X_n je náhodný výběr z rozdělení $R(\theta_1, \theta_2)$. Víme, že

$$E X_i = \frac{\theta_1 + \theta_2}{2} \quad \text{a} \quad \text{var } X_i = \frac{(\theta_2 - \theta_1)^2}{12}.$$

Odhadovací soustava rovnic v tomto případě je

$$\overline{X}_n = \frac{\widehat{\theta}_{1n} + \widehat{\theta}_{2n}}{2}, \quad \text{var } X_i = \frac{(\widehat{\theta}_{2n} - \widehat{\theta}_{1n})^2}{12}.$$

Vyřešením této soustavy dostáváme

$$\widehat{\theta}_{1n} = \overline{X}_n - \sqrt{3S_n^2} \quad \text{a} \quad \widehat{\theta}_{2n} = \overline{X}_n + \sqrt{3S_n^2}.$$

Jelikož z Věty 2.6 víme,

$$\sqrt{n} \left[\begin{pmatrix} \overline{X}_n \\ S_n^2 \end{pmatrix} - \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix} \right] \xrightarrow[n \rightarrow \infty]{d} N_2(\mathbf{0}, \Sigma),$$

kde $\Sigma = \begin{pmatrix} \sigma^2 & \sigma^3 \gamma_3 \\ \sigma^3 \gamma_3 & \sigma^4 (\gamma_4 - 1) \end{pmatrix}$ a $\gamma_3 = \frac{\mathbb{E}(X_i - \mu)^3}{\sigma^3}$, tak pomocí Δ -metody lze ukázat, že

$$\sqrt{n} \begin{bmatrix} \widehat{\theta}_{1n} \\ \widehat{\theta}_{2n} \end{bmatrix} - \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \xrightarrow{d} N_2(\mathbf{0}, \mathbb{D}\Sigma\mathbb{D}^T),$$

kde \mathbb{D} je Jakobiho matice zobrazení $\tau^{-1}(x_1, x_2) = (x_1 - \sqrt{3x_2}, x_1 + \sqrt{3x_2})$ v bodě $(\mathbb{E} X_i, \text{var } X_i)$. Tudíž odhad $\widehat{\theta}_n = (\widehat{\theta}_{1n}, \widehat{\theta}_{2n})$ je asymptoticky normální (a tedy dle tvrzení 1.4 také konsistentní).

6. X_1, \dots, X_n je náhodný výběr z bea rozdělení $B(\alpha, \beta)$ (viz např. kapitola 8.2.7 [Kulich, 2018](#)), tj. $\mathbb{E} X_i = \frac{\alpha}{\alpha+\beta}$ a $\text{var } X_i = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$. Momentovou metodou dostaneme konsistentní a asymptoticky normální odhady

$$\widehat{\alpha}_n = \bar{X}_n \left(\frac{\bar{X}_n(1 - \bar{X}_n)}{S_n^2} - 1 \right) \quad \text{a} \quad \widehat{\beta}_n = (1 - \bar{X}_n) \left(\frac{\bar{X}_n(1 - \bar{X}_n)}{S_n^2} - 1 \right)$$

(odhady jsou smysluplné pouze pokud $S_n^2 < \bar{X}_n(1 - \bar{X}_n)$).

Poznámka.

- Odhady získané momentovou metodou mívají větší asymptotický rozptyl než odhady metodou maximální věrohodnosti, která bude probírána v Matematické statistice 2.
- Pomocí věty o implicitní funkci se dá dokázat, že stačí, aby τ měla spojitou derivaci na nějakém okolí bodu $(\mathbb{E} X_i, \text{var } X_i)$.

3.4. INTERVALOVÝ ODHAD

Máme náhodný výběr $\mathbf{X} = (X_1, X_2, \dots, X_n)$, model \mathcal{F} a parametr $\theta = t(F) \in \mathbb{R}$ pro $F \in \mathcal{F}$, který chceme v daném modelu odhadnout. Nechť $F_X \in \mathcal{F}$ je skutečné rozdělení náhodného vektoru \mathbf{X} a $\theta_X \equiv t(F_X)$ je skutečná hodnota hledaného parametru.

3.4.1. DEFINICE

Definice 3.5 Interval $B_n = B_n(\mathbf{X}) \subset \mathbb{R}$ se nazývá *intervalový odhad* parametru $\theta_X \in \mathbb{R}$ o *spolehlivosti* $1 - \alpha$ v modelu \mathcal{F} , právě když

$$P[\omega \in \Omega : B_n(\mathbf{X}(\omega)) \ni \theta_X] = 1 - \alpha, \quad \text{pro každé rozdělení } F_X \in \mathcal{F}.$$

Interval B_n se nazývá *asymptotický intervalový odhad* parametru $\theta_X \in \mathbb{R}$ o (*přibližné*) *spolehlivosti* $1 - \alpha$ v modelu \mathcal{F} , právě když

$$P[\omega \in \Omega : B_n(\mathbf{X}(\omega)) \ni \theta_X] \xrightarrow{n \rightarrow \infty} 1 - \alpha \quad \text{pro každé rozdělení } F_X \in \mathcal{F}.$$

Poznámka.

- Interval B_n je náhodný (spočítaný z dat), zatímco parametr θ_X je pevný. Výraz $B_n \ni \theta_X$ čteme „interval B_n pokrývá (skutečnou hodnotu) θ_X “.
- Intervalovému odhadu se běžně říká i jinak, např. *interval spolehlivosti s pravděpodobností pokrytí (s koeficientem spolehlivosti) $1-\alpha$ nebo $(1-\alpha)100$ -procentní konfidenční interval* pro parametr θ_X .^{*} Číslo $\alpha \in (0, 1)$ je předem zvolené; obvykle se bere $\alpha = 0,05$ a počítají se 95procentní intervaly. Můžeme se však setkat i s intervaly, jež mají pokrytí 90 % či 99 %.
- Ne vždy je možné či vhodné počítat přesné intervaly spolehlivosti. Často se spokojujeme s intervaly asymptotickými, jejichž pokrytí se pro velké rozsahy výběru blíží k požadované hodnotě.
- Intervalové odhady zde definujeme pouze pro reálné parametry. Podobný koncept však lze zavést i pro vektorové parametry; hledáme náhodnou množinu B_n , která pokrývá skutečnou hodnotu se zadanou pravděpodobností. Této množině pak říkáme *oblast spolehlivosti*[†]. Tvar množiny B_n lze ale potom volit mnoha různými způsoby.

Poznámka. Rozeznáváme intervalové odhady oboustranné a jednostranné (levo- a pravostranné).

- Interval tvaru $(\eta_L(\mathbf{X}), \eta_U(\mathbf{X}))$, kde $\eta_L(\mathbf{X})$ a $\eta_U(\mathbf{X})$ jsou dvě náhodné veličiny splňující $P[\eta_L(\mathbf{X}) < \eta_U(\mathbf{X})] = 1$, $\eta_L(\mathbf{X}) > -\infty$ a $\eta_U(\mathbf{X}) < \infty$ s.j., nazýváme *oboustranný interval spolehlivosti*. Obvykle jej sestavujeme tak, aby platilo (alespoň asymptoticky)

$$P[\theta_X \leq \eta_L(\mathbf{X})] = \frac{\alpha}{2}, \quad P[\theta_X \geq \eta_U(\mathbf{X})] = \frac{\alpha}{2}.$$

- Interval tvaru $(\eta_L(\mathbf{X}), \infty)$ nazýváme *levostranný (dolní) interval spolehlivosti*. Máme $P[\eta_L(\mathbf{X}) < \theta_X] = 1 - \alpha$.
- Interval tvaru $(-\infty, \eta_U(\mathbf{X}))$ nazýváme *pravostranný (horní) interval spolehlivosti*. Máme $P[\theta_X < \eta_U(\mathbf{X})] = 1 - \alpha$.

Příklad (střední hodnota normálního rozdělení se známým rozptylem). Vezměme si problém intervalového odhadu střední hodnoty pro normálně rozdělená data se známým rozptylem.

Data: $X_1, \dots, X_n \sim F_X$

Model: $F_X \in \mathcal{F} = \{N(\mu, \sigma_X^2), \mu \in \mathbb{R}, \sigma_X^2 \text{ známo}\}$

Odhadovaný parametr: $\theta_X = E X_i \equiv \mu_X$

Postup:

1. Máme bodový odhad \bar{X}_n , který je nestranný a konsistentní pro μ_X . Víme, že $\bar{X}_n \sim N(\mu_X, \sigma_X^2/n)$. Tudiž

$$\frac{\sqrt{n}(\bar{X}_n - \mu_X)}{\sigma_X} \sim N(0, 1).$$

^{*} Angl. *confidence interval with coverage probability/confidence level $1 - \alpha$* [†] Angl. *confidence set*

2. Vyjdeme z rovnosti

$$P\left[u_{\frac{\alpha}{2}} < \frac{\sqrt{n}(\bar{X}_n - \mu_X)}{\sigma_X} < u_{1-\alpha/2}\right] = 1 - \alpha,$$

kde $u_\alpha = \Phi^{-1}(\alpha)$ je α -kvantil normovaného normálního rozdělení, a postupnými úpravami (s využitím symetrie hustoty $N(0, 1)$ kolem 0) dojdeme k

$$P\left[\bar{X}_n - u_{1-\alpha/2} \frac{\sigma_X}{\sqrt{n}} < \mu_X < \bar{X}_n + u_{1-\alpha/2} \frac{\sigma_X}{\sqrt{n}}\right] = 1 - \alpha.$$

3. Získali jsme oboustranný interval spolehlivosti (η_L, η_U) . Jeho krajní body jsou

$$\eta_L(\mathbf{X}) = \bar{X}_n - u_{1-\alpha/2} \frac{\sigma_X}{\sqrt{n}}, \quad \eta_U(\mathbf{X}) = \bar{X}_n + u_{1-\alpha/2} \frac{\sigma_X}{\sqrt{n}}.$$

Kvantily normovaného normálního rozdělení, které potřebujeme pro konstrukci intervalů spolehlivosti, jsou uvedeny v Tabulce 3.1.

Pro $\alpha = 0,05$ vezmeme kvantil $u_{0,975} \doteq 1,96$ a dostaneme 95% oboustranný interval spolehlivosti. To znamená, že tento interval pokrývá skutečnou střední hodnotu μ_X s pravděpodobností 0,95.

4. Jednostranný interval bychom získali drobnou modifikací kroku 2. *Levostranný interval* vyjde

$$(\eta_L(\mathbf{X}), \infty), \quad \text{kde} \quad \eta_L(\mathbf{X}) = \bar{X}_n - u_{1-\alpha} \frac{\sigma_X}{\sqrt{n}}.$$

Pravostranný interval vyjde

$$(-\infty, \eta_U(\mathbf{X})), \quad \text{kde} \quad \eta_U(\mathbf{X}) = \bar{X}_n + u_{1-\alpha} \frac{\sigma_X}{\sqrt{n}}.$$

Jednostranné intervaly se od oboustranného liší hodnotou kvantilu normálního rozdělení (používají $u_{1-\alpha}$ namísto $u_{1-\alpha/2}$). Pro 95% jednostranný interval spolehlivosti bychom vzali kvantil $u_{0,95} \doteq 1,645$.

Poznámka. Délka intervalu spolehlivosti:

- se zkracuje s rostoucím počtem pozorování n ,
- roste s rostoucím rozptylem dat σ_X^2 ,
- roste s rostoucí pravděpodobností pokrytí $1 - \alpha$.

Tabulka 3.1.: Vybrané hodnoty kvantilů normovaného normálního rozdělení.

κ	0,9	0,95	0,975	0,99	0,995
$u_\kappa = \Phi^{-1}(\kappa)$	1,282	1,645	1,960	2,326	2,576

Příklad. Nechť X_1, \dots, X_n je náhodný výběr z rozdělení $N(\mu_X, \sigma_X^2)$, rozptyl σ_X^2 známe. Kolik pozorování potřebujeme, aby délka oboustranného intervalu spolehlivosti pro střední hodnotu μ_X nepřekročila stanovenou mez $d > 0$?

Máme $2u_{1-\alpha/2} \sigma_X / \sqrt{n} \leq d$. Tudíž potřebujeme alespoň $4u_{1-\alpha/2}^2 \sigma_X^2 / d^2$ pozorování. Za povšimnutí stojí, že pokud chceme zkrátit interval spolehlivosti na polovinu, tak musíme zvětšit rozsah výběru čtyřikrát.

Lemma 3.2 (interval spolehlivosti po transformaci parametru) Je-li (η_L, η_U) (asymptotický) interval spolehlivosti pro parametr θ_X s pravděpodobností pokrytí $1 - \alpha$ a je-li ψ rostoucí spojitá reálná funkce na parametrickém prostoru $\Theta = \{t(F), F \in \mathcal{F}\} \subseteq \mathbb{R}$, pak $(\psi(\eta_L), \psi(\eta_U))$ je (asymptotický) interval spolehlivosti pro parametr $\psi(\theta_X)$ s pravděpodobností pokrytí $1 - \alpha$.

Důkaz. Z předpokladu lemmatu vyplývá, že pro přesný interval spolehlivosti platí

$$1 - \alpha = P[\eta_L(\mathbf{X}) < \theta_X < \eta_U(\mathbf{X})] = P[\psi(\eta_L(\mathbf{X})) < \psi(\theta_X) < \psi(\eta_U(\mathbf{X}))].$$

Analogicky pro asymptotické intervaly spolehlivosti. □

Příklad. Nechť X_1, \dots, X_n je náhodný výběr z Poissonova rozdělení $Po(\lambda)$. Potom dle příkladu na straně 37 víme, že

$$\sqrt{n} \left(2\sqrt{\bar{X}_n} - 2\sqrt{\lambda_X} \right) \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Z tohoto výsledku snadno odvodíme, že asymptotický interval spolehlivosti pro $\sqrt{\lambda_X}$ je

$$\left(\sqrt{\bar{X}_n} - \frac{u_{1-\alpha/2}}{2\sqrt{n}}, \sqrt{\bar{X}_n} + \frac{u_{1-\alpha/2}}{2\sqrt{n}} \right).$$

A tedy interval spolehlivosti pro λ_X je

$$\left(\left[\max \left\{ \sqrt{\bar{X}_n} - \frac{u_{1-\alpha/2}}{2\sqrt{n}}, 0 \right\} \right]^2, \left[\sqrt{\bar{X}_n} + \frac{u_{1-\alpha/2}}{2\sqrt{n}} \right]^2 \right).$$

3.4.2. KONSTRUKCE INTERVALOVÝCH ODHADŮ

Nechť $\mathbf{X} = (X_1, \dots, X_n)$, kde X_1, X_2, \dots, X_n je náhodný výběr z rozdělení $F_X \in \mathcal{F}$. Odhadujeme parametr $\theta_X = t(F_X) \in \mathbb{R}$. Popišme si stručně obecný postup při konstrukci oboustranných intervalových odhadů pro θ_X .

1. Nalezneme funkci $\varphi(x, \theta_X)$ takovou, že φ je prostá a spojitá funkce v argumentu θ_X pro každé x a rozdělení náhodné veličiny $Z_n \equiv \varphi(\mathbf{X}, \theta_X)$ je známé alespoň asymptoticky (nezávisí ani na θ_X ani na jiných neznámých parametrech) a je nedegenerované. Náhodná veličina Z_n se nazývá *pivotální*. Při konstrukci funkce φ můžeme vyjít např. z bodového odhadu parametru θ_X , jehož rozdělení většinou známe alespoň asymptoticky. Označíme F_Z (přesnou či asymptotickou) distribuční funkci Z_n a $c_\alpha = F_Z^{-1}(\alpha)$ budiž α -kvantil rozdělení F_Z .

3. Odhadování parametrů

2. Vyjdeme z rovnosti

$$P(c_{\alpha/2} < \varphi(\mathbf{X}, \theta_X) < c_{1-\alpha/2}) = 1 - \alpha \quad (\text{nebo } \rightarrow 1 - \alpha)$$

a „osamostatníme“ θ_X . Za tímto účelem potřebujeme zinvertovat $\varphi(x, \theta)$ jakožto funkci argumentu θ při pevném x . Tj. nechť existuje $\bar{\varphi}(x, t)$ taková, že

$$\varphi(x, \bar{\varphi}(x, t)) = t \quad \text{a} \quad \bar{\varphi}(x, \varphi(x, \theta)) = \theta$$

pro všechna x, t a θ . Jelikož funkce $\bar{\varphi}(x, t)$ je zpravidla klesající funkcí druhého argumentu t , tak dostáváme

$$P(\bar{\varphi}(\mathbf{X}, c_{1-\alpha/2}) < \theta_X < \bar{\varphi}(\mathbf{X}, c_{\alpha/2})) = 1 - \alpha.$$

3. Získali jsme (asymptotickou) interval spolehlivosti $(\eta_L(\mathbf{X}), \eta_U(\mathbf{X}))$ s pravděpodobností pokrytí $1 - \alpha$, kde $\eta_L(\mathbf{X}) = \bar{\varphi}(\mathbf{X}, c_{1-\alpha/2})$ a $\eta_U(\mathbf{X}) = \bar{\varphi}(\mathbf{X}, c_{\alpha/2})$.

Příklad (rozptyl a směrodatná odchylka normálního rozdělení). Vezměme si problém intervalového odhadu směrodatné odchylky v normálním rozdělení.

Data: $X_1, \dots, X_n \sim F_X$

Model: $F_X \in \mathcal{F} = \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$

Odhadovaný parametr: $\sigma_X = \sqrt{\text{var } \bar{X}_i}$

Postup:

Zabývejme se nejprve rozptylem σ_X^2 . Jeho nestranný a konsistentní odhad je S_n^2 . Z věty 2.8, část (i), víme, že

$$\frac{(n-1)S_n^2}{\sigma_X^2} \sim \chi_{n-1}^2.$$

Vezmeme tedy $Z_n = (n-1)S_n^2/\sigma_X^2$, $F_Z = \chi_{n-1}^2$ a $c_\alpha = \chi_{n-1}^2(\alpha)$, tj. α -kvantil rozdělení χ_{n-1}^2 (viz Tabulka 3.2).

Tabulka 3.2.: Vybrané hodnoty kvantilů $\chi_f^2(\kappa)$ rozdělení χ^2 s f stupni volnosti.

f	κ							
	0,01	0,025	0,05	0,1	0,9	0,95	0,975	0,99
5	0,554	0,831	1,145	1,610	9,236	11,070	12,833	15,086
10	2,558	3,247	3,940	4,865	15,987	18,307	20,483	23,209
15	5,229	6,262	7,261	8,547	22,307	24,996	27,488	30,578
25	11,524	13,120	14,611	16,473	34,382	37,652	40,646	44,314
100	70,065	74,222	77,929	82,358	118,498	124,342	129,561	135,807

Vyjdeme z rovnosti

$$P \left[\chi_{n-1}^2(\alpha/2) < \frac{(n-1)S_n^2}{\sigma_X^2} < \chi_{n-1}^2(1-\alpha/2) \right] = 1 - \alpha$$

a postupnými úpravami dojdeme k

$$P \left[\frac{(n-1)S_n^2}{\chi_{n-1}^2(1-\alpha/2)} < \sigma_X^2 < \frac{(n-1)S_n^2}{\chi_{n-1}^2(\alpha/2)} \right] = 1 - \alpha.$$

Získali jsme interval spolehlivosti

$$\left(\frac{(n-1)S_n^2}{\chi_{n-1}^2(1-\alpha/2)}, \frac{(n-1)S_n^2}{\chi_{n-1}^2(\alpha/2)} \right) \quad (3.3)$$

pro rozptyl σ_X^2 s pravděpodobností pokrytí $1 - \alpha$.

Interval spolehlivosti pro směrodatnou odchylku σ_X získáme aplikováním odmocniny na krajní body intervalu pro rozptyl

$$\left(\frac{\sqrt{n-1} S_n}{\sqrt{\chi_{n-1}^2(1-\alpha/2)}}, \frac{\sqrt{n-1} S_n}{\sqrt{\chi_{n-1}^2(\alpha/2)}} \right),$$

viz také Lemma 3.2 (odmocnina je rostoucí a spojitá funkce na $(0, \infty)$).

Příklad (střední hodnota normálního rozdělení s neznámým rozptylem). Vezměme si problém intervalového odhadu střední hodnoty pro normálně rozdělená data s neznámým rozptylem.

Data: $X_1, \dots, X_n \sim F_X$

Model: $F_X \in \mathcal{F} = \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$

Odhadovaný parametr: $\theta_X = E X_i \equiv \mu_X$

Postup:

Odhad \bar{X}_n je nestranný a konsistentní pro μ_X , odhad S_n^2 je nestranný a konsistentní pro $\sigma_X^2 \equiv \text{var } X_i$. Z věty 2.10 víme, že

$$T_n = \frac{\sqrt{n} (\bar{X}_n - \mu_X)}{S_n} \sim t_{n-1}.$$

Vezmeme tedy T_n jako pivotální náhodnou veličinu, F_Z je distribuční funkce rozdělení t_{n-1} a $c_\alpha = t_{n-1}(\alpha)$ (α -kvantil rozdělení t_{n-1}). Vybrané kvantily t -rozdělení jsou uvedeny v Tabulce 3.3. Jak je vidět, už pro $n - 1 = 25$ jsou jen o málo větší než kvantily normovaného normálního rozdělení, k nimž konvergují při počtu stupňů volnosti rostoucím nade všechny meze. Větší hodnoty t -kvantilů proti kvantilům normovaného normálního rozdělení používaným v úvodním příkladě odrážejí zvýšenou variabilitu pivotální statistiky způsobenou neznalostí skutečného rozptylu.

3. Odhadování parametrů

Tabulka 3.3.: Vybrané hodnoty kvantilů $t_f(\kappa)$ rozdělení t s f stupni volnosti.

f	κ				
	0,9	0,95	0,975	0,99	0,995
5	1,476	2,015	2,571	3,365	4,032
10	1,372	1,812	2,228	2,764	3,169
15	1,341	1,753	2,131	2,602	2,947
25	1,316	1,708	2,060	2,485	2,787
100	1,290	1,660	1,984	2,364	2,626
∞	1,282	1,645	1,960	2,326	2,576

Vydeme z rovnosti

$$P\left[t_{n-1}\left(\frac{\alpha}{2}\right) < \frac{\sqrt{n}(\bar{X}_n - \mu_X)}{S_n} < t_{n-1}\left(1 - \frac{\alpha}{2}\right)\right] = 1 - \alpha$$

a stejným postupem jako u normálního rozdělení se známým rozptylem dojdeme k intervalu

$$\left(\bar{X}_n - t_{n-1}\left(1 - \frac{\alpha}{2}\right) \frac{S_n}{\sqrt{n}}, \bar{X}_n + t_{n-1}\left(1 - \frac{\alpha}{2}\right) \frac{S_n}{\sqrt{n}}\right), \quad (3.4)$$

který má pravděpodobnost pokrytí přesně $1 - \alpha$.

Příklad (střední hodnota libovolného rozdělení s konečným rozptylem). Vezměme si problém intervalového odhadu střední hodnoty bez předpokladu normality dat.

Data: $X_1, \dots, X_n \sim F_X$

Model: $F_X \in \mathcal{F} = \mathcal{L}_+^2$ (všechna rozdělení s konečným a nenulovým rozptylem)

Odhadovaný parametr: $\theta_X = E X_i \equiv \mu_X$

Postup:

Odhad \bar{X}_n je nestranný a konsistentní pro μ_X , odhad S_n^2 je nestranný a konsistentní pro $\sigma_X^2 \equiv \text{var } X_i$. Z věty 2.9 víme, že

$$T_n = \frac{\sqrt{n}(\bar{X}_n - \mu_X)}{S_n} \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Vezmeme tedy T_n jako pivotální statistiku.

Vydeme z limitního vztahu (zdůvodněného konvergencí v distribuci pivotální veličiny)

$$P\left[u_{\frac{\alpha}{2}} < \frac{\sqrt{n}(\bar{X}_n - \mu_X)}{S_n} < u_{1-\alpha/2}\right] \xrightarrow[n \rightarrow \infty]{} 1 - \alpha.$$

Tedy asymptotický interval spolehlivosti by byl

$$\left(\bar{X}_n - u_{1-\alpha/2} \frac{S_n}{\sqrt{n}}, \bar{X}_n + u_{1-\alpha/2} \frac{S_n}{\sqrt{n}}\right). \quad (3.5)$$

3. Odhadování parametrů

Jelikož pro $n \rightarrow \infty$ kvantil $t_{n-1}(\alpha)$ konverguje k u_α (pro libovolné $0 < \alpha < 1$), tak máme, že také interval (3.4), který byl přesným intervalem spolehlivosti pro μ_X u výběru z normálního rozdělení, je zároveň asymptotickým intervalem spolehlivosti pro μ_X pro data pocházející z jakéhokoli rozdělení s konečným nenulovým rozptylem.

Všimněme si, že $|t_{n-1}(\alpha)| > |u_\alpha|$ pro všechna $n \geq 2$, tudíž interval (3.4) je delší než interval (3.5). Z důvodu opatrnosti se tedy doporučuje používat spíše interval (3.4).

Příklad (alternativní rozdělení). Ukažme si nyní jeden možný způsob odvození asymptotického intervalového odhadu pro pravděpodobnost úspěchu v alternativním rozdělení. (Několik dalších intervalových odhadů pro tento problém si ukážeme později.)

Data: $X_1, \dots, X_n \sim F_X$

Model: $F_X \in \mathcal{F} = \{\text{Alt}(p), p \in (0, 1)\}$

Odhadovaný parametr: $p_X = E X_i = P[X_i = 1]$

Postup:

Jelikož odhadujeme pravděpodobnost, vyjdeme z empirické relativní četnosti $\hat{p}_n = \bar{X}_n$, která je nestranným a konsistentním odhadem p (věta 2.3). Z centrální limitní věty (tvrzení 1.6) víme, že $\sqrt{n}(\hat{p}_n - p_X) \xrightarrow[n \rightarrow \infty]{d} N(0, p_X(1 - p_X))$. Tudíž

$$\frac{\sqrt{n}(\hat{p}_n - p_X)}{\sqrt{p_X(1 - p_X)}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Levá strana je nelineární funkcí p_X , ale můžeme si ji zjednodušit. Z konsistence \hat{p}_n a věty o spojité transformaci (tvrzení 1.2) víme, že

$$\sqrt{\hat{p}_n(1 - \hat{p}_n)} \xrightarrow[n \rightarrow \infty]{P} \sqrt{p_X(1 - p_X)}.$$

Ze Sluckého věty (tvrzení 1.3) dostaneme

$$\frac{\sqrt{n}(\hat{p}_n - p_X)}{\sqrt{\hat{p}_n(1 - \hat{p}_n)}} = \frac{\sqrt{n}(\hat{p}_n - p_X)}{\sqrt{p_X(1 - p_X)}} \frac{\sqrt{p_X(1 - p_X)}}{\sqrt{\hat{p}_n(1 - \hat{p}_n)}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1). \quad (3.6)$$

Vezmeme tedy $Z_n = \frac{\sqrt{n}(\hat{p}_n - p_X)}{\sqrt{\hat{p}_n(1 - \hat{p}_n)}}$, $F_Z = \Phi$ a $c_\alpha = u_\alpha$ (α -kvantil normovaného normálního rozdělení).

Vyjdeme z limitního vztahu

$$P \left[-u_{1-\alpha/2} < \frac{\sqrt{n}(\hat{p}_n - p_X)}{\sqrt{\hat{p}_n(1 - \hat{p}_n)}} < u_{1-\alpha/2} \right] \xrightarrow[n \rightarrow \infty]{} 1 - \alpha$$

a postupnými úpravami dojdeme k

$$P \left[\hat{p}_n - u_{1-\alpha/2} \frac{\sqrt{\hat{p}_n(1 - \hat{p}_n)}}{\sqrt{n}} < p_X < \hat{p}_n + u_{1-\alpha/2} \frac{\sqrt{\hat{p}_n(1 - \hat{p}_n)}}{\sqrt{n}} \right] \xrightarrow[n \rightarrow \infty]{} 1 - \alpha.$$

Získali jsme tedy interval

$$\left(\hat{p}_n - u_{1-\alpha/2} \frac{\sqrt{\hat{p}_n(1-\hat{p}_n)}}{\sqrt{n}}, \hat{p}_n + u_{1-\alpha/2} \frac{\sqrt{\hat{p}_n(1-\hat{p}_n)}}{\sqrt{n}} \right),$$

jehož pravděpodobnost pokrytí konverguje k $1 - \alpha$ pro $n \rightarrow \infty$.

3.5. EMPIRICKÉ ODHADY

Mějme dán náhodný výběr X_1, X_2, \dots, X_n z rozdělení F_X . Ukažme si, jak lze odhadnout některé charakteristiky rozdělení F_X .

3.5.1. EMPIRICKÁ DISTRIBUČNÍ FUNKCE

Zabývejme se nejprve situací, kdy máme jednorozměrná pozorování X_1, \dots, X_n a odhadujeme celou distribuční funkci $F_X(x)$ pro $x \in \mathbb{R}$. Pracujeme s modelem, který zahrnuje veškerá rozdělení na \mathbb{R} , tj. na distribuční funkci F_X neklademe vůbec žádné podmínky.

Definice 3.6 Funkci $\hat{F}_n(x) \stackrel{\text{df}}{=} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq x\}$ nazýváme *empirická distribuční funkce** náhodného výběru X_1, X_2, \dots, X_n .

Poznámka. Hodnota \hat{F}_n v bodě x je rovna počtu pozorování, která nepřekročí x , dělenému celkovým počtem pozorování. Funkce \hat{F}_n je neklesající, zprava spojitá, po částech konstantní, skáče v pozorovaných hodnotách veličin X_i , velikosti skoků jsou dány počtem pozorování rovných x děleným celkovým počtem pozorování. Empirická distribuční funkce má všechny vlastnosti distribuční funkce diskrétního rozdělení.

Pro pevné x je hodnota $\hat{F}_n(x)$ vlastně relativní četnost jevu $[X_i \leq x]$ spočítaná z n pozorování, přičemž pravděpodobnost tohoto jevu je $F_X(x)$. Z věty 2.3 rovnou dostaneme nejdůležitější vlastnosti empirické distribuční funkce.

Věta 3.3 (vlastnosti empirické distribuční funkce) Pro libovolné $x \in \mathbb{R}$ platí:

- (i) $E \hat{F}_n(x) = F_X(x)$ (nestrannost), $\text{var}(\hat{F}_n(x)) = \frac{F_X(x)[1-F_X(x)]}{n}$;
- (ii) $\hat{F}_n(x) \xrightarrow[n \rightarrow \infty]{P} F_X(x)$ (bodová konsistence);
- (iii) $\sqrt{n} [\hat{F}_n(x) - F_X(x)] \xrightarrow[n \rightarrow \infty]{d} N(0, F_X(x)[1 - F_X(x)])$ (asymptotická normalita);
- (iv) $n\hat{F}_n(x) \sim \text{Bi}(n, F_X(x))$;
- (v) $\sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F_X(x)| \xrightarrow[n \rightarrow \infty]{P} 0$ (stejněměrná konsistence).

Poznámka.

* Angl. *empirical distribution function*

- Z bodu (iii) předchozí věty lze odvodit asymptotický interval spolehlivosti pro $F_X(x)$ stejně jako v případě parametru alternativního rozdělení (viz str. 56).
- Bod (v) se někdy nazývá Glivenkova-Cantelliho věta. Nelze jej odvodit z věty 2.3 ani jiných výsledků, které máme k dispozici. Bude dokázán na jedné z pokročilejších přednášek z teorie pravděpodobnosti.

3.5.2. IDEA EMPIRICKÝCH ODHADŮ

Z empirické distribuční funkce lze odvodit odhady mnoha základních charakteristik rozdělení F_X . Nechť $\theta_X = t(F_X)$ je hledaný parametr. Umíme-li jej spočítat ze skutečné distribuční funkce F_X , můžeme jej stejným způsobem spočítat i z empirické distribuční funkce \widehat{F}_n . Dostaneme tak odhad $\widehat{\theta}_n \stackrel{\text{df}}{=} t(\widehat{F}_n)$. Těmto odhadům říkáme *empirické odhady*. Uvidíme, že v řadě případů mají empirické odhady rozumné vlastnosti.

Ukažme si tento postup nejprve na příkladě empirického odhadu střední hodnoty. Máme

$$E X_i = \int_{-\infty}^{\infty} x dF_X(x).$$

Empirický odhad střední hodnoty získáme dosazením \widehat{F}_n na místo neznámé funkce F_X . Dostaneme

$$\int_{-\infty}^{\infty} x d\widehat{F}_n(x) = \int_{-\infty}^{\infty} x d\left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq x\}\right) = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} x d\mathbb{1}\{X_i \leq x\} = \frac{1}{n} \sum_{i=1}^n X_i,$$

kde jsme využili toho, že $G(x) = \mathbb{1}\{X_i \leq x\}$ je pro pevné X_i vlastně distribuční funkcí konstanty nabývající hodnoty X_i s pravděpodobností 1. Došli jsme tedy k tomu, že empirickým odhadem střední hodnoty je aritmetický průměr, o němž již víme, že je nestranný a konsistentní.

Poznámka. Zafixujme si $\omega \in \Omega$ a označme si realizace náhodných veličin jako $x_1 = X_1(\omega), \dots, x_n = X_n(\omega)$. Pak \widehat{F}_n splňuje vlastnosti distribuční funkce. Pokud nyní Y je nějaká náhodná veličina s distribuční funkcí \widehat{F}_n , pak integrál $\int_{-\infty}^{\infty} x d\widehat{F}_n(x)$ je vlastně střední hodnota Y . Jelikož rozdělení dané distribuční funkcí \widehat{F}_n je diskrétní a splňuje $P(Y = x_i) = \frac{1}{n}$ pro všechna $i = 1, \dots, n$, tak

$$E Y = \sum_{i=1}^n x_i P(Y = x_i) = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{i=1}^n X_i(\omega).$$

3.5.3. EMPIRICKÉ ODHADY MOMENTŮ

Nechť X_1, X_2, \dots, X_n je náhodný výběr z rozdělení F_X a h je měřitelná reálná funkce taková, že $E |h(X_i)| < \infty$. Dá se snadno ověřit, že empirickým odhadem parametru $E h(X_i)$ je průměr naměřených hodnot $h(X_i)$, tj. $\frac{1}{n} \sum_{i=1}^n h(X_i)$. Tento odhad je nestranný a konsistentní.

3. Odhadování parametrů

Odvodme si *empirický odhad rozptylu* $\sigma_X^2 = E X_i^2 - (E X_i)^2$. Víme, že empirickým odhadem $E X_i$ je \bar{X}_n a empirickým odhadem $E X_i^2$ je $\frac{1}{n} \sum_{i=1}^n X_i^2$. Empirický odhad rozptylu tedy je

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

Poznámka. Platí

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{n}{n-1} \hat{\sigma}_n^2.$$

Pro velká n je rozdíl mezi $\hat{\sigma}_n^2$ a S_n^2 malý, neboť s pomocí věty 2.6(i)

$$\hat{\sigma}_n^2 - S_n^2 = -\frac{S_n^2}{n} \xrightarrow[n \rightarrow \infty]{P} 0.$$

Jak plyne z věty 2.6, výběrový rozptyl S_n^2 je nestranný a konsistentní odhad σ_X^2 . Empirický odhad rozptylu $\hat{\sigma}_n^2$ je konsistentní, ale není nestranný. Na druhou stranu z příkladu na straně 42 víme, že v modelu $\mathcal{F} = \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$ platí $MSE(\hat{\sigma}_n^2) < MSE(S_n^2)$.

Podobně můžeme odvodit empirické odhady pro momenty vyšších řádů. *Empirické odhady necentrálních momentů* $\mu'_k = E X_i^k$ jsou

$$\hat{\mu}'_k = \frac{1}{n} \sum_{i=1}^n X_i^k.$$

Empirické odhady centrálních momentů $\mu_k = E (X_i - E X_i)^k$ jsou

$$\hat{\mu}_k = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^k.$$

Empirické necentrální momenty jsou evidentně nestranné a konsistentní. Empirické centrální momenty jsou konsistentní, nikoli však obecně nestranné.

Empirický odhad šikmosti je

$$\hat{\gamma}_3 = \frac{\hat{\mu}_3}{(\hat{\sigma}_n^2)^{3/2}},$$

empirický odhad špičatosti je

$$\hat{\gamma}_4 = \frac{\hat{\mu}_4}{\hat{\sigma}_n^4}.$$

Oba jsou konsistentní (z věty o spojitě transformaci, tvrzení 1.2).

Cvičení. Dokažte, že pokud $E |X_i|^k < \infty$, pak $\hat{\mu}_k \xrightarrow[n \rightarrow \infty]{P} \mu_k$.

Návod:

$$\hat{\mu}_k = \frac{1}{n} \sum_{i=1}^n \sum_{j=0}^k \binom{j}{k} X_i^k (-\bar{X}_n)^{k-j} = \sum_{j=0}^k \binom{j}{k} \left(\frac{1}{n} \sum_{i=1}^n X_i^k \right) (-\bar{X}_n)^{k-j}.$$

3.5.4. EMPIRICKÉ ODHADY KVANTILŮ

Nechť α je předem dané číslo z intervalu $(0, 1)$. *Kvantilová funkce* rozdělení F_X je definována jako

$$F_X^{-1}(\alpha) = \inf \{x : F_X(x) \geq \alpha\}.$$

Pak α -kvantilem rozdělení F_X rozumíme číslo $u_X(\alpha) = F_X^{-1}(\alpha)$. Pro α -kvantil platí

$$F_X(u_X(\alpha)) \geq \alpha \quad \text{a} \quad F_X(u_X(\alpha) - h) < \alpha \quad \text{pro } \forall h > 0.$$

Jako empirický odhad použijeme hodnotu α -kvantilu empirické distribuční funkce, tedy

$$\widehat{F}_n^{-1}(\alpha) = \inf \{x : \widehat{F}_n(x) \geq \alpha\}.$$

Definice 3.7 (Výběrový kvantil) Pro $\alpha \in (0, 1)$ definujeme *empirický (výběrový) α -kvantil** jako $\widehat{u}_n(\alpha) = \widehat{F}_n^{-1}(\alpha)$.

Poznámka.

- Připomeňme si, že empirická distribuční funkce je po částech konstantní se skoky v bodech $X_{(1)}, X_{(2)}, \dots, X_{(n)}$. Tudíž empirický kvantil dle naší definice bude vhodně vybraná pořádková statistika. Jelikož dále platí,

$$\widehat{F}_n(X_{(k)}) \geq \frac{k}{n} \quad \text{a} \quad \widehat{F}_n(X_{(k)} - h) < \frac{k}{n} \quad \text{pro } \forall h > 0,$$

tak empirický kvantil bude splňovat

$$\widehat{u}_n(\alpha) = X_{(k_\alpha)}, \quad \text{kde} \quad k_\alpha = \begin{cases} n\alpha & \text{pro } (n\alpha) \in \mathbb{N}, \\ \lfloor n\alpha \rfloor + 1 & \text{pro } (n\alpha) \notin \mathbb{N}. \end{cases}$$

Jelikož nepředpokládáme spojitost rozdělení, tak pořádkové statistice $X_{(k_\alpha)}$ je třeba rozumět ve smyslu poznámky na straně 35.

- Pro $\alpha = 0,5$ dostaneme *výběrový medián*[†]: $\widehat{m}_n = X_{(\frac{n+1}{2})}$ pro n liché a $\widehat{m}_n = X_{(\frac{n}{2})}$ pro n sudé.
- Výběrový α -kvantil splňuje nerovnosti

$$\widehat{F}_n(\widehat{u}_n(\alpha)) \geq \alpha \quad \text{a} \quad \lim_{h \searrow 0} \widehat{F}_n(\widehat{u}_n(\alpha) - h) < \alpha,$$

tj. alespoň $n\alpha$ pozorování je menší nebo rovno $\widehat{u}_n(\alpha)$ a zároveň pro všechna $h > 0$ je alespoň $n(1 - \alpha)$ pozorování větší nebo rovno $\widehat{u}_n(\alpha) - h$.

- Existuje mnoho různých definic výběrového α -kvantilu (zpravidla jako nějaké lineární interpolace mezi body $X_{(k_\alpha-1)}, X_{(k_\alpha)}$ a $X_{(k_\alpha+1)}$). Např. pro sudá n se výběrový medián často definuje jako

$$\widehat{m}_n = \frac{X_{(\frac{n}{2})} + X_{(\frac{n}{2}+1)}}{2}.$$

* Angl. *empirical quantile, sample quantile* † Angl. *sample median*

3. Odhadování parametrů

Následující lemma charakterizuje výběrový kvantil jako řešení minimalizačního problému. Ve srovnání s lemmatem 2.1 však tato charakterizace není nutně jednoznačná (viz diskuse níže).

Lemma 3.4 Necht' $\alpha \in (0, 1)$. Pro výběrový α -kvantil $\hat{u}_n(\alpha)$ platí

$$\hat{u}_n(\alpha) = \arg \min_{c \in \mathbb{R}} \sum_{i=1}^n \varrho_\alpha(X_i - c),$$

kde $\varrho_\alpha(u) = \alpha u \mathbb{1}\{u \geq 0\} + (1 - \alpha)(-u) \mathbb{1}\{u < 0\}$.

Všimněme si, že pro $\alpha = \frac{1}{2}$ dostáváme $\varrho_{1/2}(u) = \frac{1}{2}|u|$. Jelikož konstanta $\frac{1}{2}$ je pro optimalizační úlohu nepodstatná, tak pro výběrový medián platí

$$\hat{m}_n = \arg \min_{c \in \mathbb{R}} \sum_{i=1}^n |X_i - c|,$$

tj. \hat{m}_n minimalizuje součet absolutních odchylek. Nicméně lze snadno nahlédnout, že výše uvedená charakterizace mediánu nemusí být jednoznačná. Např. pro $n = 2$ je součet $\sum_{i=1}^2 |X_i - c|$ minimalizován pro jakékoliv $c \in \langle X_{(1)}, X_{(2)} \rangle$.

Obecně pak lze dokázat, že pokud $k_\alpha = (n\alpha) \in \mathbb{N}$, pak jakékoliv c z uzavřeného intervalu $\langle X_{(k_\alpha)}, X_{(k_\alpha+1)} \rangle$ minimalizuje $\sum_{i=1}^n \varrho_\alpha(X_i - c)$.

Poznámka. Minimalizační problém z části (ii) lze psát jako úlohu lineárního programování ve tvaru

$$\arg \min_{c \in \mathbb{R}} \left[-(1 - \alpha) \sum_{i: X_i < c} (X_i - c) + \alpha \sum_{i: X_i \geq c} (X_i - c) \right].$$

Zavedeme-li značení $U_i = (X_i - c) \mathbb{1}\{X_i \geq c\}$, $V_i = -(X_i - c) \mathbb{1}\{X_i < c\}$, $\mathbf{U} = (U_1, \dots, U_n)^\top$, $\mathbf{V} = (V_1, \dots, V_n)^\top$, $\mathbf{X} = (X_1, \dots, X_n)^\top$, můžeme problém přepsat jako úlohu lineárního programování ve $(2n + 1)$ -dimensionálním prostoru

$$\min_{\mathbf{U}, \mathbf{V}, c} \alpha \mathbf{1}_n^\top \mathbf{U} + (1 - \alpha) \mathbf{1}_n^\top \mathbf{V}$$

při omezeních

$$c \mathbf{1}_n + \mathbf{U} - \mathbf{V} = \mathbf{X}, \quad \mathbf{U} \geq 0, \quad \mathbf{V} \geq 0.$$

Tento minimalizační problém samozřejmě nemusí mít právě jedno řešení. Minima může být dosaženo na celém intervalu hodnot.

Vlastnosti výběrového kvantilu budeme dokazovat pouze pro spojitá rozdělení s rostoucí distribuční funkcí F_X a hustotou f_X .

Věta 3.5 Necht' $\alpha \in (0, 1)$. Necht' X_1, \dots, X_n je náhodný výběr z rozdělení, která má distribuční funkcí F_X spojitou a rostoucí na nějakém okolí bodu $u_X(\alpha)$.

3. Odhadování parametrů

(i) Potom $\widehat{u}_n(\alpha) \xrightarrow[n \rightarrow \infty]{P} u_X(\alpha)$.

(ii) Pokud navíc existuje hustota f_X , která je spojitá a nenulová v bodě $u_X(\alpha)$, pak

$$\sqrt{n}[\widehat{u}_n(\alpha) - u_X(\alpha)] \xrightarrow[n \rightarrow \infty]{d} N(0, V(\alpha)), \quad \text{kde} \quad V(\alpha) = \frac{\alpha(1-\alpha)}{f_X^2(u_X(\alpha))}.$$

Důkaz. Část (i): Necht' $\varepsilon > 0$. Potřebujeme ukázat, že

$$P(|\widehat{u}_n(\alpha) - u_X(\alpha)| > \varepsilon) \xrightarrow[n \rightarrow \infty]{} 0.$$

K tomu nám stačí ukázat, že

$$P(\widehat{u}_n(\alpha) < u_X(\alpha) - \varepsilon) \xrightarrow[n \rightarrow \infty]{} 0 \quad \text{a zároveň} \quad P(\widehat{u}_n(\alpha) > u_X(\alpha) + \varepsilon) \xrightarrow[n \rightarrow \infty]{} 0.$$

Počítejme tedy

$$\begin{aligned} P(\widehat{u}_n(\alpha) < u_X(\alpha) - \varepsilon) &= P(X_{(k_\alpha)} < u_X(\alpha) - \varepsilon) \\ &= P\left(\sum_{i=1}^n \mathbb{1}\{X_i < u_X(\alpha) - \varepsilon\} \geq k_\alpha\right) \\ &\leq P\left(\widehat{F}_n(u_X(\alpha) - \varepsilon) - F_X(u_X(\alpha) - \varepsilon) \geq \frac{k_\alpha}{n} - F_X(u_X(\alpha) - \varepsilon)\right). \end{aligned} \quad (3.7)$$

Z věty 3.3 nyní plyne, že

$$\widehat{F}_n(u_X(\alpha) - \varepsilon) - F_X(u_X(\alpha) - \varepsilon) \xrightarrow[n \rightarrow \infty]{P} 0, \quad (3.8)$$

a zároveň z předpokladů dokazované věty

$$\frac{k_\alpha}{n} - F_X(u_X(\alpha) - \varepsilon) \xrightarrow[n \rightarrow \infty]{} \alpha - F_X(u_X(\alpha) - \varepsilon) > 0. \quad (3.9)$$

Kombinací (3.8) a (3.9) pak dostáváme, že pravá strana rovnosti (3.7) konverguje k nule, tudíž jsme dokázali, že $P(\widehat{u}_n(\alpha) < u_X(\alpha) - \varepsilon) \xrightarrow[n \rightarrow \infty]{} 0$.

Podobně

$$\begin{aligned} P(\widehat{u}_n(\alpha) > u_X(\alpha) + \varepsilon) &= P\left(\sum_{i=1}^n \mathbb{1}\{X_i \leq u_X(\alpha) + \varepsilon\} < k_\alpha\right) \\ &\leq P\left(\widehat{F}_n(u_X(\alpha) + \varepsilon) - F_X(u_X(\alpha) + \varepsilon) < \frac{k_\alpha}{n} - F_X(u_X(\alpha) + \varepsilon)\right). \end{aligned} \quad (3.10)$$

Z věty 3.3 nyní plyne, že

$$\widehat{F}_n(u_X(\alpha) + \varepsilon) - F_X(u_X(\alpha) + \varepsilon) \xrightarrow[n \rightarrow \infty]{P} 0, \quad (3.11)$$

a zároveň z předpokladů dokazované věty

$$\frac{k_\alpha}{n} - F_X(u_X(\alpha) + \varepsilon) \xrightarrow[n \rightarrow \infty]{} \alpha - F_X(u_X(\alpha) + \varepsilon) < 0. \quad (3.12)$$

3. Odhadování parametrů

Kombinací (3.11) a (3.12) pak dostáváme, že pravá strana rovnosti (3.10) konverguje k nule, tudíž jsme dokázali, že $P(\widehat{u}_n(\alpha) > u_X(\alpha) + \varepsilon) \xrightarrow{n \rightarrow \infty} 0$.

Část (ii): Podobně jako v části (i) počítejme

$$\begin{aligned} P\left(\sqrt{n}[\widehat{u}_n(\alpha) - u_X(\alpha)] \leq x\right) &= P\left(\widehat{u}_n(\alpha) \leq u_X(\alpha) + \frac{x}{\sqrt{n}}\right) \\ &= P\left(\widehat{F}_n\left(u_X(\alpha) + \frac{x}{\sqrt{n}}\right) - F_X\left(u_X(\alpha) + \frac{x}{\sqrt{n}}\right) \geq \frac{k_\alpha}{n} - F_X\left(u_X(\alpha) + \frac{x}{\sqrt{n}}\right)\right) \\ &= P(Z_n \geq x_n), \end{aligned}$$

kde

$$Z_n = \frac{\sqrt{n}\left[\widehat{F}_n\left(u_X(\alpha) + \frac{x}{\sqrt{n}}\right) - F_X\left(u_X(\alpha) + \frac{x}{\sqrt{n}}\right)\right]}{\sqrt{\alpha(1-\alpha)}}$$

a

$$x_n = \frac{\sqrt{n}\left[\frac{k_\alpha}{n} - F_X\left(u_X(\alpha) - \frac{x}{\sqrt{n}}\right)\right]}{\sqrt{\alpha(1-\alpha)}}.$$

Z centrální limitní věty pro trojúhelníková schéma (např. Věta 16.4 Lachout, 2004) pak plyne, že $Z_n \xrightarrow[n \rightarrow \infty]{d} Z$, kde $Z \sim N(0, 1)$. Dále z předpokladů věty dostáváme $x_n \xrightarrow[n \rightarrow \infty]{} \frac{-x f_X(u_X(\alpha))}{\sqrt{\alpha(1-\alpha)}}$. Tedy celkem máme

$$P\left(\sqrt{n}[\widehat{u}_n(\alpha) - u_X(\alpha)] \leq x\right) \xrightarrow{n \rightarrow \infty} P\left(Z \geq \frac{-x f_X(u_X(\alpha))}{\sqrt{\alpha(1-\alpha)}}\right) = P\left(Z \leq \frac{x f_X(u_X(\alpha))}{\sqrt{\alpha(1-\alpha)}}\right),$$

což společně s definicí konvergence v distribuci implikuje tvrzení věty. □

Příklad. Uvažujte náhodný výběr X_1, \dots, X_n z rozdělení s hustotou

$$f_X(x) = \frac{1}{2} \mathbb{1}\{x \in (0, 1) \cup (2, 3)\}.$$

Ukažte, že v tomto případě není výběrový medián konzistentním odhadem mediánu $m = u_X\left(\frac{1}{2}\right) = 1$. Který z předpokladů věty 3.5(i) není splněn?

Ukažte, že $P(\widehat{m}_n \geq 2) \rightarrow \frac{1}{2}$ pro $n \rightarrow \infty$.

Intervalový odhad pro $u_X(\alpha)$

Asymptotický rozptyl $V(\alpha)$ výběrového kvantilu (viz věta 3.5(ii)) se špatně odhaduje, protože nemáme k dispozici univerzálně použitelný a spolehlivý odhad hustoty. Za předpokladu, že F_X je spojitá v $u_X(\alpha)$ lze ke konstrukci intervalu spolehlivosti využít pořádkové statistiky.

3. Odhadování parametrů

Např. *oboustranný interval spolehlivosti* pro $u_X(\alpha)$ s pravděpodobností pokrytí $1-\beta$ hledáme ve tvaru $(X_{(k_L)}, X_{(k_U)})$. Pro určení čísel k_L a k_U si všimněme, že

$$P(X_{(k_L)} \geq u_X(\alpha)) = P\left(\sum_{i=1}^n \mathbb{1}\{X_i < u_X(\alpha)\} \leq k_L - 1\right) = P(\text{Bi}(n, \alpha) \leq k_L - 1),$$

$$P(X_{(k_U)} \leq u_X(\alpha)) = P\left(\sum_{i=1}^n \mathbb{1}\{X_i \leq u_X(\alpha)\} \geq k_U\right) = P(\text{Bi}(n, \alpha) \geq k_U).$$

Tedy čísla k_L a k_U můžeme určit pomocí binomického rozdělení jako nejvyšší a nejmenší možné přirozené číslo takové, aby

$$P(\text{Bi}(n, \alpha) \leq k_L - 1) \leq \frac{\beta}{2}, \quad P(\text{Bi}(n, \alpha) \geq k_U) \leq \frac{\beta}{2}.$$

V případě, že nemáme možnost pracovat přímo s binomickým rozdělením, tak můžeme využít normální aproximaci binomického rozdělení. V tomto případě je dobré si povšimnout, že

$$P(\text{Bi}(n, \alpha) \leq k_L - 1) = P(\text{Bi}(n, \alpha) < k_L) \quad \text{a} \quad P(\text{Bi}(n, \alpha) \geq k_U) = P(\text{Bi}(n, \alpha) > k_U - 1).$$

Proto jako „kompromis“ před normální aproximací vycházíme z rovností

$$P(X_{(k_L)} \geq u_X(\alpha)) = P(\text{Bi}(n, \alpha) < k_L - \frac{1}{2}), \quad P(X_{(k_U)} \leq u_X(\alpha)) = P(\text{Bi}(n, \alpha) > k_U - \frac{1}{2}).$$

Nyní pomocí normální aproximace

$$P(\text{Bi}(n, \alpha) < k_L - \frac{1}{2}) = P\left(\frac{\text{Bi}(n, \alpha) - n\alpha}{\sqrt{n\alpha(1-\alpha)}} < \frac{k_L - \frac{1}{2} - n\alpha}{\sqrt{n\alpha(1-\alpha)}}\right) \doteq \Phi\left(\frac{k_L - \frac{1}{2} - n\alpha}{\sqrt{n\alpha(1-\alpha)}}\right),$$

$$P(\text{Bi}(n, \alpha) > k_U - \frac{1}{2}) = P\left(\frac{\text{Bi}(n, \alpha) - n\alpha}{\sqrt{n\alpha(1-\alpha)}} > \frac{k_U - \frac{1}{2} - n\alpha}{\sqrt{n\alpha(1-\alpha)}}\right) \doteq 1 - \Phi\left(\frac{k_U - \frac{1}{2} - n\alpha}{\sqrt{n\alpha(1-\alpha)}}\right).$$

Odtud pak již vyjádříme přibližné hodnoty k_L a k_U

$$k_L = \left\lceil \frac{1}{2} + n\alpha - u_{1-\frac{\beta}{2}} \sqrt{n\alpha(1-\alpha)} \right\rceil, \quad k_U = \left\lceil \frac{1}{2} + n\alpha + u_{1-\frac{\beta}{2}} \sqrt{n\alpha(1-\alpha)} \right\rceil.$$

„Kompromis“ popsany výše se zpravidla nazývá *oprava na spojitost**. Tato „oprava“ však nespočívá v tom, že bychom něco dělali spojitým, ale je to jistá opatrnost v případě, že diskrétní rozdělení (v tomto případě binomické) aproximuje spojitým rozdělením (v tomto případě normálním).

Poznámka. Může se stát, že pro malé rozsahy výběrů n a α blízké nule nebo jedné je jedna z pravděpodobností $P(\text{Bi}(n, \alpha) = 0) > \frac{\beta}{2}$ resp. $P(\text{Bi}(n, \alpha) = n) > \frac{\beta}{2}$. V takovém případě volíme za dolní (resp. horní) mez intervalu spolehlivosti $-\infty$ (resp. $+\infty$).

Cvičení. Ukažte, že pokud bychom vynechali předpoklad spojitosti distribuční funkce v odhadovaném kvantilu $u_X(\alpha)$, tak uzavřený interval $\langle X_{(k_L)}, X_{(k_U)} \rangle$ bude mít (pro dostatečně velké n) pravděpodobnost pokrytí alespoň $1 - \beta$.

* Angl. *continuity correction*

3.5.5. EMPIRICKÉ ODHADY PRO NÁHODNÉ VEKTORY

Empirické odhady prvních dvou momentů můžeme snadno rozšířit na náhodné vektory. Nechť $\mathbf{X}_1, \dots, \mathbf{X}_n$ je náhodný výběr nezávislých k -rozměrných náhodných vektorů s rozdělením F_X . Jednotlivé složky vektoru \mathbf{X}_i budeme značit X_{ij} , $i = 1, \dots, n$, $j \in \{1, \dots, k\}$. Dále označme

$$\boldsymbol{\mu} = E \mathbf{X}_i, \quad \Sigma = \text{var } \mathbf{X}_i.$$

Empirickým odhadem $\boldsymbol{\mu}$ je zřejmě vektor empirických odhadů jeho jednotlivých složek, čili k -rozměrný výběrový průměr

$$\bar{\mathbf{X}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i.$$

Empirický odhad rozptylové matice Σ bychom dostali z vyjádření

$$\Sigma = E (\mathbf{X}_i - E \mathbf{X}_i)(\mathbf{X}_i - E \mathbf{X}_i)^\top = E \mathbf{X}_i \mathbf{X}_i^\top - (E \mathbf{X}_i)(E \mathbf{X}_i)^\top = E \mathbf{X}_i^{\otimes 2} - (E \mathbf{X}_i)^{\otimes 2}$$

a nahrazením středních hodnot jejich empirickými odhady (tj. průměry) bychom dostali

$$\hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^{\otimes 2} - \bar{\mathbf{X}}_n^{\otimes 2} = \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}}_n)(\mathbf{X}_i - \bar{\mathbf{X}}_n)^\top.$$

Většinou se však pracuje s tzv. *výběrovou rozptylovou maticí**, která se definuje jako vícerozměrná obdoba výběrového rozptylu S_n^2 :

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}}_n)(\mathbf{X}_i - \bar{\mathbf{X}}_n)^\top.$$

Poznámka.

- S_n^2 má na diagonále výběrové rozptyly jednotlivých složek, tj.

$$S_j^2 = \frac{1}{n-1} \sum_{i=1}^n (X_{ij} - \bar{X}_j)^2,$$

pro $j \in \{1, \dots, k\}$, kde $\bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{ij}$.

- Prvek (j, m) matice S_n^2 je dán výrazem

$$S_{jm} = \frac{1}{n-1} \sum_{i=1}^n (X_{ij} - \bar{X}_j)(X_{im} - \bar{X}_m)$$

pro $j \in \{1, \dots, k\}$ a $m \in \{1, \dots, k\}$, $j \neq m$. Tato náhodná veličina odhaduje kovarianci $\text{cov}(X_{ij}, X_{im})$ mezi j -tou a m -tou složkou \mathbf{X}_i . Říkáme jí *výběrová kovariance*.

* Angl. *sample covariation matrix*

- \mathbb{S}_n^2 je pozitivně semidefinitní a platí

$$\mathbb{S}_n^2 = \frac{n}{n-1} \widehat{\Sigma}_n = \frac{n}{n-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^{\otimes 2} - \overline{\mathbf{X}}_n^{\otimes 2} \right).$$

Následující tvrzení ukazuje, že jak $\overline{\mathbf{X}}_n$ tak \mathbb{S}_n^2 jsou nestranné a konsistentní odhady.

Tvrzení 3.6

(i) Je-li $E |X_{ij}| < \infty$ pro všechna $j \in \{1, \dots, k\}$, pak $E \overline{\mathbf{X}}_n = \boldsymbol{\mu}$ a $\overline{\mathbf{X}}_n \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\mu}$.

(ii) Je-li $\text{var}(X_{ij}) < \infty$ pro všechna $j \in \{1, \dots, k\}$, pak $E \mathbb{S}_n^2 = \Sigma$ a $\mathbb{S}_n^2 \xrightarrow[n \rightarrow \infty]{P} \Sigma$.

Důkaz. Část (i): Plyne přímo z věty 2.2 aplikované po složkách.

Část (ii): Konsistence \mathbb{S}_n^2 se ukáže analogicky jako u S_n^2 (viz věta 2.6(i)).

Nestrannost lze dokázat např. následujícím způsobem:

$$\begin{aligned} E \mathbb{S}_n^2 &= \frac{n}{n-1} \left[\frac{1}{n} \sum_{i=1}^n E \mathbf{X}_i^{\otimes 2} - E \left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right)^{\otimes 2} \right] \\ &= \frac{n}{n-1} \left(E \mathbf{X}_i^{\otimes 2} - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n E \mathbf{X}_i \mathbf{X}_j^T \right) \\ &= \frac{n}{n-1} \left(E \mathbf{X}_i^{\otimes 2} - \frac{1}{n^2} \sum_{i=1}^n E \mathbf{X}_i^{\otimes 2} - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1, j \neq i}^n E \mathbf{X}_i \mathbf{X}_j^T \right) \\ &= \frac{n}{n-1} \left[E \mathbf{X}_i^{\otimes 2} \left(1 - \frac{1}{n} \right) - \frac{n-1}{n} (E \mathbf{X}_i)^{\otimes 2} \right] = \Sigma. \end{aligned}$$

□

Vzpomeňme si na definici korelačního koeficientu mezi veličinami X_{ij} a X_{im} :

$$\rho(X_{ij}, X_{im}) = \frac{\text{cov}(X_{ij}, X_{im})}{\sqrt{\text{var } X_{ij} \text{ var } X_{im}}}.$$

Je logické zavést výběrový korelační koeficient jakožto empirický odhad tohoto parametru vzniklý z empirických odhadů jeho jednotlivých komponent.

Definice 3.8 Výběrový korelační koeficient* $\widehat{\rho}_{jm}$ veličin X_{ij} a X_{im} , $j \in \{1, \dots, k\}$ a $m \in \{1, \dots, k\}$, $j \neq m$, definujeme jako

$$\widehat{\rho}_{jm} = \frac{S_{jm}}{S_j S_m} = \frac{\sum_{i=1}^n (X_{ij} - \overline{X}_j)(X_{im} - \overline{X}_m)}{\sqrt{\sum_{i=1}^n (X_{ij} - \overline{X}_j)^2 \sum_{i=1}^n (X_{im} - \overline{X}_m)^2}}.$$

* Angl. *sample correlation coefficient*

Poznámka.

- $-1 \leq \hat{\varrho}_{jm} \leq 1$ (viz Cauchyho-Schwarzova nerovnost).
- $\hat{\varrho}_{jm} = 1$ (resp. -1) právě když existují konstanty $a \in \mathbb{R}$ a $b > 0$ (resp. $b < 0$) takové, že $X_{ij} = a + bX_{im}$ pro všechna $i = 1, \dots, n$.
- $\hat{\varrho}_{jm}$ je konsistentní odhad korelačního koeficientu $\varrho(X_{ij}, X_{im})$ (to plyne z konzistence \mathbb{S}_n^2 a tvrzení 1.2), ale není nestranný.

Cvičení. Podrobně dokažte, že $\hat{\varrho}_{jm} \xrightarrow[n \rightarrow \infty]{P} \varrho(X_{ij}, X_{im})$.

Přípravné příklady ke zkoušce.

1. Mějme náhodný výběr X_1, \dots, X_n z rozdělení s hustotou $f(x; \delta) = \frac{e^{-x/\delta}}{\delta} \mathbb{1}\{x > 0\}$, kde $\delta > 0$ je neznámý parametr. Uvažujte odhad $\widehat{\delta}_n = \overline{X}_n$. Ukažte, že tento odhad je nestraný odhad parametru δ_X . Dále uvažujte odhad $\widetilde{\delta}_n(a) = a \overline{X}_n$, kde a je konstanta. Najděte a tak, aby minimalizovalo střední čtvercovou chybu odhadu $\widetilde{\delta}_n(a)$.
2. Mějme náhodný výběr X_1, \dots, X_n z alternativního rozdělení s parametrem p_X . Odhadněte parametr p_X momentovou metodou a transformací tohoto odhadu vytvořte odhad parametru $\theta_X = p_X(1 - p_X)$. Prozkoumejte nestrannost a konsistenci takto vytvořeného odhadu rozptylu. Jak se tento odhad liší od obyčejného výběrového rozptylu?
3. Mějme náhodný výběr X_1, \dots, X_n z alternativního rozdělení s parametrem p_X . Potom z příkladu na straně 56 víme, že asymptotický intervalový odhad o spolehlivosti $1 - \alpha$ pro parametr p_X je

$$\left(\widehat{p}_n - u_{1-\alpha/2} \frac{\sqrt{\widehat{p}_n(1-\widehat{p}_n)}}{\sqrt{n}}, \widehat{p}_n + u_{1-\alpha/2} \frac{\sqrt{\widehat{p}_n(1-\widehat{p}_n)}}{\sqrt{n}} \right).$$

Odvoďte odsud intervalový odhad pro $\theta_X = p_X(1 - p_X)$.

Nechť intervalový odhad o spolehlivosti $1 - \alpha$ pro parametr p_X spočtený z daných data vyšel (0.35, 0.55). Jak by v tomto případě vypadal intervalový odhad pro $\theta_X = p_X(1 - p_X)$?

4. Nechť X_1, \dots, X_n je náhodný výběr z rozdělení $N(\mu_X, 9)$. Jaké musí být n , aby délka intervalu spolehlivosti pro μ_X s pravděpodobností pokrytí 0,90 byla nejvýše 0,25?
5. Nechť X_1, \dots, X_n je náhodný výběr z rozdělení $N(\mu_X, \sigma_X^2)$, kde σ_X^2 je neznámé. Sestavte za této situace interval spolehlivosti pro μ_X . Vyjádřete vzorec pro pravděpodobnost, že pro dané n a dané σ_X^2 bude šířka tohoto intervalu spolehlivosti s pravděpodobností pokrytí $1 - \alpha$ menší než 10.
6. Nechť \overline{X}_n je průměr náhodného výběru X_1, \dots, X_n z rozdělení $Po(\lambda_X)$. Určete asymptotické rozdělení odhadu výběrového průměru \overline{X}_n a na základě tohoto rozdělení určete asymptotický interval spolehlivosti pro parametr $\theta_X = \exp\{-\lambda_X\}$.
7. Nechť X_1, \dots, X_n je náhodný výběr z rovnoměrného rozdělení $R(0, 1)$. Nechť $k_n = \lceil n/2 - \sqrt{n} \rceil$. Dokažte, že $X_{(k_n)} \xrightarrow[n \rightarrow \infty]{P} 1/2$.

4. PRINCIPY TESTOVÁNÍ HYPOTÉZ

4.1. ZÁKLADNÍ POJMY A DEFINICE

Nechť X_1, \dots, X_n je náhodný výběr nezávislých k -rozměrných náhodných vektorů s rozdělením $F_X \in \mathcal{F}$, kde \mathcal{F} je model. Nechť $\theta = t(F) \in \mathbb{R}^d$ je charakteristika rozdělení, která nás zajímá (parametr), nechť $\Theta = \{t(F), F \in \mathcal{F}\} \subseteq \mathbb{R}^d$ označuje všechny možné hodnoty parametru v modelu \mathcal{F} (nazývá se *parametrický prostor*^{*}). Označme skutečný parametr jako $\theta_X = t(F_X)$. Označme celá napozorovaná data symbolem $X = (X_1, \dots, X_n)$.

Příklady. Nově zaváděné pojmy a tvrzení budeme v celé této kapitole objasňovat na následujících příkladech.

- A. Nechť X_1, \dots, X_n je náhodný výběr z rozdělení $N(\theta_X, \sigma_0^2)$, kde $\sigma_0^2 > 0$ je známo. Máme tedy model

$$\mathcal{F}^A = \{N(\theta, \sigma_0^2), \theta \in \mathbb{R}\}.$$

- B. Nechť X_1, \dots, X_n je náhodný výběr z rozdělení $N(\theta_X, \sigma_X^2)$, kde σ_X^2 není známo. Pracujeme s modelem

$$\mathcal{F}^B = \{N(\theta, \sigma^2), \theta \in \mathbb{R}, \sigma^2 > 0\} \supset \mathcal{F}^A.$$

- C. Nechť X_1, \dots, X_n je náhodný výběr z rozdělení F_X s konečným a nenulovým rozptylem. Pracujeme s neparametrickým modelem

$$\mathcal{F}^C = \mathcal{L}_+^2 \supset \mathcal{F}^B \supset \mathcal{F}^A.$$

Testovaným parametrem bude střední hodnota $\theta = \int x dF(x)$, jeho skutečná hodnota je $\theta_X = E X_i$, dimenze d parametru θ je 1. Parametrický prostor je $\Theta = \mathbb{R}$.

Zvolme si nyní dvě neprázdné disjunktní podmnožiny Θ , které označíme Θ_0 a Θ_1 . Řekněme, že nás nyní nezajímá konkrétní hodnota parametru θ_X , ale chceme pouze odpovědět na otázku, zdali $\theta_X \in \Theta_0$ nebo $\theta_X \in \Theta_1$.

Definice 4.1 (Hypotéza a alternativa)

- Množinu Θ_0 nazýváme [nulová] *hypotéza*[†], množinu Θ_1 nazýváme *alternativa*[‡] (nebo také alternativní hypotéza).

* Angl. *parameter space* † Angl. *null hypothesis* ‡ Angl. *alternative hypothesis*

- Označme

$$\mathcal{F}_0 \stackrel{\text{df}}{=} \{F \in \mathcal{F} : t(F) \in \Theta_0\},$$

tj. všechna rozdělení v modelu \mathcal{F} , jejichž parametry splňují hypotézu. Jestliže $\mathcal{F}_0 = \{F_0\}$ (tj. v modelu existuje právě jedno rozdělení, které hypotézu splňuje), hypotézu nazýváme *jednoduchou*^{*}, jinak *složenou*[†].

- Označme

$$\mathcal{F}_1 \stackrel{\text{df}}{=} \{F \in \mathcal{F} : t(F) \in \Theta_1\},$$

tj. všechna rozdělení v modelu \mathcal{F} , jejichž parametry splňují alternativu. Jestliže $\mathcal{F}_1 = \{F_1\}$ (tj. v modelu existuje právě jedno rozdělení, které alternativu splňuje), alternativu nazýváme *jednoduchou*[‡], jinak *složenou*[§].

Poznámka.

- Hypotézu označujeme obvykle symbolem H_0 , alternativu symbolem H_1 . Mluvíme o *testování hypotézy*

$$H_0 : \theta_X \in \Theta_0 \quad \text{proti alternativě} \quad H_1 : \theta_X \in \Theta_1.$$

- Jednoduchou hypotézu tedy dostaneme, pokud $\Theta_0 = \{\theta_0\}$ je jednobodová množina a zároveň existuje právě jedno rozdělení $F_0 \in \mathcal{F}$ takové, že $t(F_0) = \theta_0$.
- Jednoduchou alternativu tedy dostaneme, pokud $\Theta_1 = \{\theta_1\}$ je jednobodová množina a zároveň existuje právě jedno rozdělení $F_1 \in \mathcal{F}$ takové, že $t(F_1) = \theta_1$.

Většinou bereme $\Theta_1 = \Theta_0^c$ a $\mathcal{F}_1 = \mathcal{F}_0^c$. Pokud tomu tak není, tj. $\Theta_0 \cup \Theta_1 \subsetneq \Theta$, tak si můžeme model zúžit na $\mathcal{F}^0 = \{F \in \mathcal{F} : t(F) \in \Theta_0 \cup \Theta_1\}$. Předpokládat, že $\Theta_1 = \Theta_0^c$ a $\mathcal{F}_1 = \mathcal{F}_0^c$ tedy není na újmu obecnosti.

Volba hypotéz pro jednorozměrný parametr θ

- Nejobvyklejší volba hypotézy je $\Theta_0 = \{\theta_0\}$ pro nějaké předem zvolené $\theta_0 \in \mathbb{R}$, tj. testujeme $H_0 : \theta_X = \theta_0$. Za alternativu volíme $\Theta_1 = \Theta_0^c$, tj. $H_1 : \theta_X \neq \theta_0$. Výslednou proceduru pak nazýváme *oboustranný test*[¶], respektive *test proti oboustranné alternativě*.
- Jiná možnost je volit $\Theta_0 = (-\infty, \theta_0)$, tj. testovat $H_0 : \theta_X \leq \theta_0$ proti $H_1 : \theta_X > \theta_0$, případně $\Theta_0 = \langle \theta_0, \infty)$, tj. testovat $H_0 : \theta_X \geq \theta_0$ proti $H_1 : \theta_X < \theta_0$. Tyto testy nazýváme *jednostranné testy*^{||}, respektive *testy proti jednostranné alternativě*. Všimněte si, že **krajní hodnota θ_0 je pokaždé zahrnuta v nulové hypotéze**.

Volba hypotézy je dána podstatou praktického problému, který řešíme. V některých případech volíme hypotézu značně odlišně od tří zmíněných možností. V této přednášce se však budeme zabývat pouze výše zmíněnými oboustrannými a jednostrannými testy.

Příklady. Uvažujme oboustranný test parametru $\theta = t(F) = \int x dF(x) \in \mathbb{R}$. Testujeme hypotézu $H_0 : \theta_X = \theta_0$ proti alternativě $H_1 : \theta_X \neq \theta_0$.

* Angl. *simple null hypothesis* † Angl. *composite null hypothesis* ‡ Angl. *simple alternative* § Angl. *composite alternative* ¶ Angl. *two-sided test* || Angl. *one-sided tests*

- A. Model $\mathcal{F}^A = \{N(\theta, \sigma_0^2), \theta \in \mathbb{R}\}$. V tomto modelu je $\mathcal{F}_0 = \{N(\theta_0, \sigma_0^2)\}$, jedná se tedy o test jednoduché hypotézy. Alternativa je složená, $\mathcal{F}_1 = \{N(\theta, \sigma_0^2), \theta \in \mathbb{R} \setminus \{\theta_0\}\}$.
- B. Model $\mathcal{F}^B = \{N(\theta, \sigma^2), \theta \in \mathbb{R}, \sigma^2 > 0\}$. V tomto modelu je hypotéza složená, $\mathcal{F}_0 = \{N(\theta_0, \sigma^2), \sigma^2 > 0\}$, alternativa je také složená, $\mathcal{F}_1 = \{N(\theta, \sigma^2), \theta \in \mathbb{R} \setminus \{\theta_0\}, \sigma^2 > 0\}$.
- C. Model $\mathcal{F}^C = \mathcal{L}_+^2$. V tomto modelu je hypotéza složená, $\mathcal{F}_0 = \{F \in \mathcal{L}_+^2 : t(F) = \theta_0\}$, alternativa je také složená, $\mathcal{F}_1 = \{F \in \mathcal{L}_+^2 : t(F) \neq \theta_0\}$.

Na základě náhodného výběru X_1, \dots, X_n chceme rozhodnout, zda H_0 platí nebo nikoli. Použijeme k tomu nějakou vhodně zvolenou funkci dat $S_n(\mathbf{X})$, které říkáme *testová statistika*^{*}, a množinu C , které říkáme *kritický obor*[†]. Testová statistika je obvykle jednorozměrná; kritický obor je pak nějaká podmnožina \mathbb{R} . Rozhodujeme se podle toho, jestli testová statistika padne do kritického oboru, či nikoli.

- Pokud $S_n(\mathbf{X}) \in C$, učiníme závěr, že *zamítáme* hypotézu H_0 ve prospěch alternativy H_1 .
- Pokud $S_n(\mathbf{X}) \notin C$, učiníme závěr, že hypotézu H_0 *nemůžeme zamítnout* ve prospěch alternativy H_1 .

Poznámka. Někteří autoři definují kritický obor jako podmnožinu výběrového prostoru, tj. v našem značení jako $S_n^{-1}(C)$. Zamítají pak hypotézu H_0 , pokud $\mathbf{X} \in S_n^{-1}(C)$.

Definice 4.2 (Test) *Statistický test* je definován pomocí testové statistiky $S_n(\mathbf{X})$, kritického oboru C a výše uvedeného pravidla pro zamítání hypotézy. Dva testy $(S_n(\mathbf{X}), C)$ a $(S_n^*(\mathbf{X}), C^*)$ nazveme *ekvivalentní* právě když $S_n(\mathbf{X}) \in C \Leftrightarrow S_n^*(\mathbf{X}) \in C^*$ skoro jistě, tj. oba testy vydávají s pravděpodobností 1 totéž rozhodnutí.

4.2. HLADINA A SÍLA TESTU

Při testování hypotéz mohou nastat čtyři situace v závislosti na tom, zdali hypotéza ve skutečnosti platí a zdali ji test zamítne.

- **Hypotéza platí, test ji nezamítne**, tj. $\theta_X \in \Theta_0$ a $S_n(\mathbf{X}) \notin C$. V tomto případě test rozhodl správně.
- **Hypotéza platí, test ji zamítne**, tj. $\theta_X \in \Theta_0$ a $S_n(\mathbf{X}) \in C$. V tomto případě test rozhodl nesprávně.
- **Hypotéza neplatí, test ji nezamítne**, tj. $\theta_X \notin \Theta_0$ a $S_n(\mathbf{X}) \notin C$. V tomto případě test rozhodl nesprávně.
- **Hypotéza neplatí, test ji zamítne**, tj. $\theta_X \notin \Theta_0$ a $S_n(\mathbf{X}) \in C$. V tomto případě test rozhodl správně.

^{*} Angl. *test statistic* [†] Angl. *critical region*

Definice 4.3 (Chyba I. a II. druhu)

- (i) Jestliže test zamítl platnou hypotézu, říkáme, že nastala *chyba I. druhu*^{*}.
- (ii) Jestliže test nezamítl neplatnou hypotézu, říkáme, že nastala *chyba II. druhu*[†].

Možné situace jsou zachyceny v tabulce 4.1.

Tabulka 4.1.: Situace, které mohou nastat při testování hypotéz.

	Nezamítáme H_0	Zamítáme H_0
H_0 platí	OK	chyba I. druhu
H_0 neplatí	chyba II. druhu	OK

Chybám I. a II. druhu se obecně nelze vyhnout. Klasický statistický přístup k testování hypotéz spočívá v tom, že **kontrolujeme pravděpodobnost chyby I. druhu**.

Co se týká chyby II. druhu, tak ideální by bylo vybrat takový test, který minimalizuje pravděpodobnost chyby II. druhu. Jelikož však pravděpodobnost chyby II. druhu závisí na zvolené alternativě, tak takovéto ideální testy existují pouze v případech, kdy alternativa není příliš velká.

4.2.1. HLADINA TESTU

Pro $F \in \mathcal{F}$ si označme

$$P_F[S_n(\mathbf{X}) \in B] = \int \mathbb{1}\{S_n(\mathbf{x}) \in B\} dF(\mathbf{x}_1) \cdots dF(\mathbf{x}_n).$$

V případě, že v modelu \mathcal{F} je jednoznačný vztah mezi parametrem $\theta \in \Theta$ a rozdělením $F \in \mathcal{F}$ pak můžeme psát

$$P_\theta[S_n(\mathbf{X}) \in B] = \int \mathbb{1}\{S_n(\mathbf{x}) \in B\} dF(\mathbf{x}_1) \cdots dF(\mathbf{x}_n), \quad (4.1)$$

kde F je rozdělení splňující $t(F) = \theta$.

Všimněme si, že (4.1) můžeme také psát v případě, když rozdělení náhodné veličiny $S_n(\mathbf{X})$ je stejné, ať již zvolíme jakékoliv F , které splňuje, že $t(F) = \theta$.

Definice 4.4 (Hladina testu) Nechť $\alpha \in (0, 1)$ je předem stanovené číslo.

- (i) Jestliže kritický obor C splňuje podmínku

$$\sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \in C] = \alpha,$$

říkáme, že test $(S_n(\mathbf{X}), C)$ má *hladinu významnosti*[‡] přesně α .

* Angl. *type I error* † Angl. *type II error* ‡ Angl. *significance level*

(ii) Jestliže kritický obor C splňuje podmínku

$$\sup_{F \in \mathcal{F}_0} \lim_{n \rightarrow \infty} P_F[S_n(\mathbf{X}) \in C] = \alpha,$$

pak říkáme, že test $(S_n(\mathbf{X}), C)$ má hladinu α *asymptoticky*.

Poznámka.

- Je-li množina $\mathcal{F}_0 = \{F_0\}$ jednobodová, pak můžeme přesnou hladinu testu psát jednodušeji

$$\alpha = P_{\theta_0}[S_n(\mathbf{X}) \in C], \quad \text{kde } \theta_0 = t(F_0).$$

- Zhruba řečeno, hladina testu je pravděpodobnost chyby prvního druhu, to jest pravděpodobnost zamítnutí platné hypotézy. Pokud hypotéza zahrnuje více než jednu hodnotu parametru, pak jde o největší možnou pravděpodobnost chyby prvního druhu.
- Test, který požadované hladiny α dosahuje přesně, budeme nazývat *přesný test*. Test, který požadované hladiny α dosahuje jen asymptoticky, budeme nazývat *asymptotický test*.
- Někteří autoři od asymptotického testu požadují splnění podmínky

$$\sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \in C] \rightarrow \alpha \quad \text{pro } n \rightarrow \infty.$$

Tato podmínka by však vyloučila, abychom např. *t*-test (viz strana 100) nazývali asymptotickým testem pro $\mathcal{F} = \mathcal{L}_+^2$.

Klasický přístup k testování hypotéz můžeme shrnout takto:

1. Předem stanovíme požadovanou hladinu testu α , kterou má test dosáhnout buď přesně nebo asymptoticky.
2. Najdeme vhodnou testovou statistiku $S_n(\mathbf{X})$.
3. Kritický obor $C = C(\alpha)$ zvolíme v závislosti na α tak, aby hladina testu (přesná nebo asymptotická) byla právě α a přitom pravděpodobnost chyby II. druhu byla co nejmenší.

Poznámka.

- Hladina testu se volí malá, v praxi se obvykle bere $\alpha = 0,05$.
- Má-li testová statistika $S_n(\mathbf{X})$ *diskrétní rozdělení*, pak není možné dosáhnout zcela libovolné hladiny α . V případě, že předepsaná hladina α je nedosažitelná, tak se spokojujeme s hladinou $\alpha' < \alpha$, která je nejbližší k původně požadovanému α . To nám zaručí, že pravděpodobnost zamítnutí platné hypotézy nemůže být větší než zvolená tolerance α .

Terminologie.

- Testu, jehož skutečná hladina je menší než požadované α , se říká *konservativní*. Testu, jehož skutečná hladina je větší než požadované α , se říká *antikonservativní*.

4.2.2. SÍLA TESTU

Definice 4.5 (Silofunkce a síla testu) Funkce

$$\beta_n(F) = P_F[S_n(\mathbf{X}) \in C]$$

zobrazující \mathcal{F} do $\langle 0, 1 \rangle$ se nazývá *silofunkce** testu.

Pokud $F \in \mathcal{F}_1$, pak číslo $\beta_n(F)$ se nazývá *síla†* testu proti alternativě F .

Poznámka.

- Síla testu je *pravděpodobnost zamítnutí neplatné hypotézy* při dané konkrétní alternativě F . Síla závisí na alternativě, pro níž ji vyhodnocujeme. Síla je rovna doplňku pravděpodobnosti chyby II. druhu do jedničky. Síla testu nemá netriviální dolní hranici; o pravděpodobnosti chyby II. druhu nemůžeme předpokládat, že je malá.
- Má-li test přesnou, resp. asymptotickou hladinu α , pak musí platit

$$\sup_{F \in \mathcal{F}_0} \beta_n(F) = \alpha, \quad \text{resp.} \quad \sup_{F \in \mathcal{F}_0} \lim_{n \rightarrow \infty} \beta_n(F) = \alpha.$$

- Pokud existuje jednoznačný vztah mezi $\theta \in \Theta$ a $F \in \mathcal{F}$ pak se zpravidla silofunkce definuje jako zobrazení parametrického prostoru Θ do $\langle 0, 1 \rangle$ dané předpisem

$$\beta_n(\theta) = P_\theta[S_n(\mathbf{X}) \in C].$$

Poznámka (Interpretace výsledku testu).

- Skončí-li test *zamítnutím hypotézy* H_0 , znamená to, že rozdělení dat neodpovídá rozdělení, jaké by data měla za platnosti hypotézy. Pravděpodobnost chybného zamítnutí v případě, že hypotéza platí, je omezena shora hladinou α , která je malá. Hypotézu H_0 vyvracíme, prokázali jsme platnost alternativy H_1 .
- Skončí-li test tím, že *hypotézu* H_0 *nemůžeme zamítnout*, znamená to, že rozdělení dat není dostatečně odlišné od rozdělení, jaké by data měla za platnosti hypotézy. Proto nemůžeme usoudit, že hypotéza H_0 platí a alternativa neplatí. Pravděpodobnost chybného rozhodnutí v případě, že hypotéza neplatí, může být značně velká. Tento výsledek tedy neznamená potvrzení platnosti hypotézy.
- Hypotéza H_0 a alternativa H_1 při testování nevystupují symetricky. Hypotézu můžeme vyvrátit ve prospěch alternativy, ale nemůžeme ji potvrdit nebo prokázat.

Abychom mohli stanovit kritický obor $C(\alpha)$, který dodržuje požadovanou hladinu α , musíme být schopni spočítat přesné nebo asymptotické rozdělení testové statistiky za platnosti hypotézy, a to nesmí záviset na neznámých charakteristikách rozdělení F_X .

Testovou statistiku $S_n(\mathbf{X})$ tedy volíme tak, aby

- (i) její rozdělení bylo *citlivé* na skutečnou hodnotu testovaného parametru θ_X ;

* Angl. *power function* † Angl. *power*

(ii) za platnosti nulové hypotézy* její rozdělení (alespoň asymptoticky) *nezáviselo na neznámých parametrech* a bylo známo (alespoň asymptoticky).

Máme-li testovou statistiku, **kritický obor** $C(\alpha)$ volíme tak, aby

- (i) byla *dodržena požadovaná hladina* testu α ;
- (ii) v kritickém oboru byly zahrnuty ty hodnoty testové statistiky, které jsou za platnosti hypotézy *méně pravděpodobné* než za platnosti alternativy.

Příklad (A1). OBOUSTRANNÝ TEST STŘEDNÍ HODNOTY NORMÁLNÍHO ROZDĚLENÍ SE ZNÁMÝM ROZPTYLEM.

Máme náhodný výběr X_1, \dots, X_n z rozdělení $F_X = N(\theta_X, \sigma_0^2) \in \mathcal{F}^A = \{N(\theta, \sigma_0^2), \theta \in \mathbb{R}\}$. Testujeme

$$H_0 : \theta_X = \theta_0 \quad \text{proti} \quad H_1 : \theta_X \neq \theta_0.$$

Testovou statistiku založíme na bodovém odhadu parametru θ_X , tj. průměru. Víme, že

$$U_n = \frac{\sqrt{n} (\bar{X}_n - \theta_0)}{\sigma_0}$$

má za platnosti hypotézy H_0 rozdělení $N(0, 1)$. Jestliže hypotéza neplatí, tj. $\theta_X - \theta_0 = \delta \neq 0$, pak

$$U_n = \frac{\sqrt{n} (\bar{X}_n - \theta_X + \theta_X - \theta_0)}{\sigma_0} = \frac{\sqrt{n} (\bar{X}_n - \theta_X)}{\sigma_0} + \frac{\sqrt{n} \delta}{\sigma_0}$$

má rozdělení $N(\nu_n, 1)$, kde $\nu_n = \frac{\sqrt{n} \delta}{\sigma_0}$. Je-li porušena hypotéza, pak se rozdělení testové statistiky posouvá pryč od nuly, a to tím dále, čím větší je n a $|\theta_X - \theta_0|$. Hodnoty testové statistiky, které jsou daleko od nuly, tedy povedou k zamítnutí hypotézy.

Kritický obor $C(\alpha)$ bude mít tvar

$$(-\infty, c_L(\alpha)) \cup (c_U(\alpha), \infty).$$

Kritické hodnoty $c_L(\alpha)$ a $c_U(\alpha)$ určíme tak, aby

$$P_{\theta_0} [U_n \in (-\infty, c_L(\alpha))] = P_{\theta_0} [U_n \in (c_U(\alpha), \infty)] = \frac{\alpha}{2}.$$

To zaručí, že hladina testu je přesně rovna α . Odtud máme díky symetrii hustoty normovaného normálního rozdělení $c_U(\alpha) = -c_L(\alpha) = u_{1-\alpha/2}$. Test tedy funguje takto

$$\text{zamítne } H_0 : \theta_X = \theta_0 \iff |U_n| = \left| \frac{\sqrt{n} (\bar{X}_n - \theta_0)}{\sigma_0} \right| \geq u_{1-\alpha/2},$$

tj. zamítáme hypotézu, pokud se \bar{X}_n liší od hypotetické hodnoty θ_0 o více než $\frac{u_{1-\alpha/2} \sigma_0}{\sqrt{n}}$. Za kvantil $u_{1-\alpha/2}$ dosazujeme 1,96 pro $\alpha = 0,05$ a 1,645 pro $\alpha = 0,1$. Kritický obor a hustoty testové statistiky za hypotézy a za alternativy jsou zobrazeny na obrázku 4.1.

* V případě jednostranných testů bychom měli přesněji říkat, pokud skutečná hodnota parametru je na hranici nulové hypotézy a alternativy.

4. Principy testování hypotéz

Spočítejme nyní *silofunkci* tohoto testu. Vezměme nějaké θ takové, že $\theta - \theta_0 = \delta \neq 0$. Pokud θ je skutečný parametr, pak rozdělení U_n je $N(v_n, 1)$ a rozdělení $U_n - v_n$ je $N(0, 1)$. Dostaneme tedy

$$\begin{aligned} \beta_n(\theta) &= P_\theta[U_n \in C(\alpha)] = P_\theta[U_n \leq -u_{1-\alpha/2}] + P_\theta[U_n \geq u_{1-\alpha/2}] = \\ &= P_\theta[U_n - v_n \leq -u_{1-\alpha/2} - v_n] + P_\theta[U_n - v_n \geq u_{1-\alpha/2} - v_n] = \\ &= \Phi(-u_{1-\alpha/2} - v_n) + 1 - \Phi(u_{1-\alpha/2} - v_n). \end{aligned}$$

Protože $\Phi(-x) = 1 - \Phi(x)$, tento výsledek můžeme přepsat do tvaru

$$\beta_n(\theta) = \Phi(-u_{1-\alpha/2} - |v_n|) + 1 - \Phi(u_{1-\alpha/2} - |v_n|). \quad (4.2)$$

Pro $\theta = \theta_0$ dostaneme $v_n = 0$, a tedy $\beta_n(\theta_0) = \alpha$. Průběh silofunkce tohoto testu je zakreslen na obrázku 4.2.

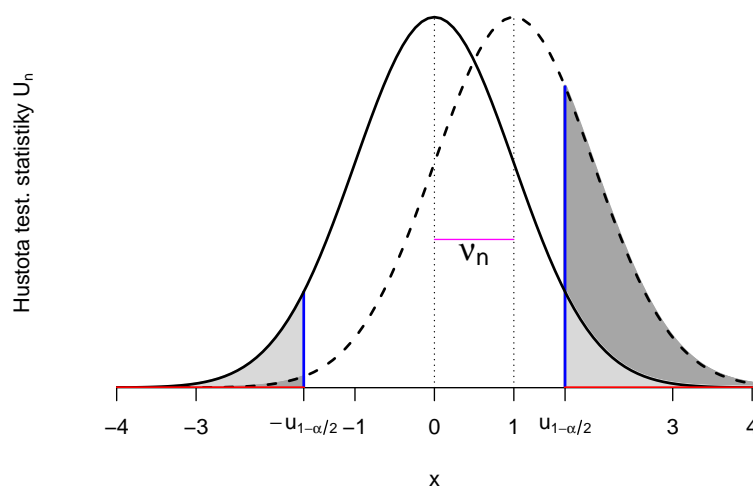
Nechť δ je nenulové. Pak $|v_n|$ roste do nekonečna s rostoucím n a ukazuje se, že od určitého n je $\Phi(-u_{1-\alpha/2} - |v_n|)$ ve srovnání s $\Phi(u_{1-\alpha/2} - |v_n|)$ zanedbatelné. Silofunkci tedy můžeme aproximovat výrazem $1 - \Phi(u_{1-\alpha/2} - \frac{\sqrt{n}|\delta|}{\sigma_0})$, přičemž platí

$$\beta_n(\theta) \geq 1 - \Phi\left(u_{1-\alpha/2} - \frac{\sqrt{n}|\delta|}{\sigma_0}\right). \quad (4.3)$$

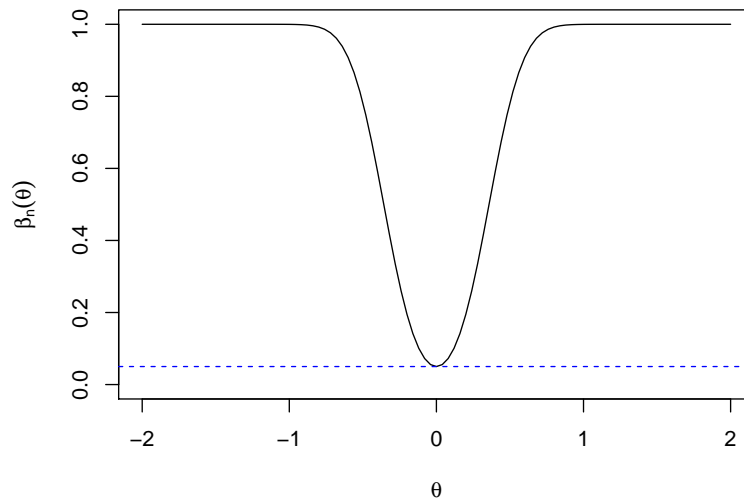
Vyřešením rovnice

$$1 - \Phi\left(u_{1-\alpha/2} - \frac{\sqrt{n}|\delta|}{\sigma_0}\right) \stackrel{!}{=} \beta,$$

Obrázek 4.1.: Hustota testové statistiky U_n za hypotézy a za alternativy pro $v_n = 1$ a $\alpha = 0,1$. Kritické hodnoty jsou vyznačeny modře, kritický obor červeně.



Obrázek 4.2.: Silofunkce oboustranného testu střední hodnoty normálního rozdělení se známým rozptylem pro $\theta_0 = 0$, $\sigma_0^2 = 1$, $n = 30$ a $\alpha = 0,05$.



pak můžeme snadno spočítat, kolik pozorování je potřeba, aby test dosáhl síly alespoň β (například 0,8). Požadovaný rozsah výběru je

$$n \geq (u_{1-\alpha/2} - u_{1-\beta})^2 \frac{\sigma_0^2}{\delta^2} = (u_{1-\alpha/2} + u_\beta)^2 \frac{\sigma_0^2}{\delta^2}. \quad (4.4)$$

Poznámka. Jak jsme viděli v předchozím příkladě, síla testu závisí na

- hladině testu α ;
- alternativě θ , respektive její vzdálenosti δ od hypotézy θ_0 ;
- rozptylu pozorování σ_0^2 ;
- počtu pozorování n .

Z těchto faktorů je možné ovlivnit pouze počet pozorování. Chceme-li dosáhnout dostatečné síly, musíme získat alespoň takový počet pozorování, jaký je uveden v (4.4).

Poznámka. Všimněme si, že síla předchozího testu proti libovolné alternativě konverguje k 1 při $n \rightarrow \infty$ (viz (4.3)). Tuto vlastnost nazýváme *konsistence testu*. Konsistence je velmi žádoucí vlastnost, jinak totiž nemusíme být schopni dosáhnout požadované síly ani při velmi velkém počtu pozorování.

Definice 4.6 Test $(S_n(\mathbf{X}), C)$ na hladině α nazveme *konsistentním testem**, jestliže $\forall F \in \mathcal{F}_1$ platí $\lim_{n \rightarrow \infty} \beta_n(F) = 1$.

Zaveďme ještě jednu užitečnou vlastnost testů: *nestrannost*.

* Angl. *consistent test*

Definice 4.7 Test $(S_n(\mathbf{X}), C)$ na hladině α nazveme *nestranným testem**, jestliže $\forall F \in \mathcal{F}_1$ platí $\beta_n(F) \geq \alpha$.

Poznámka.

- Nenechte se zmást: pojmy nestrannost a konsistence testu mají jen velmi volný (pokud vůbec nějaký) vztah k pojům nestrannost a konsistence odhadu.
- Nestrannost testu vyžaduje, aby síla proti každé alternativě byla alespoň α . Kdyby tomu tak nebylo, t.j. $\exists F \in \mathcal{F}_1$ taková, že $\beta_n(F) < \alpha$, test by tuto F vlastně považoval za součást hypotézy.
- Test, který vždy zamítá H_0 s pravděpodobností α (bez ohledu na data) je nestranný. Nestranný test tedy existuje.
- Někdy se pojmy konsistence a nestrannost vztahují vůči specifickým alternativám. Tedy například říkáme, že daný test je *konsistentní vůči konkrétní $F \in \mathcal{F}_1$* , jestliže platí $\lim_{n \rightarrow \infty} \beta_n(F) = 1$.

4.2.3. VOLBA KRITICKÉHO OBORU

Kritický obor $C(\alpha)$ má ve většině případů jeden z následujících tvarů:[†]

- $\langle c_U(\alpha), \infty \rangle$, tj. zamítáme pro *příliš velké* hodnoty testové statistiky $S_n(\mathbf{X})$;
- $(-\infty, c_L(\alpha))$, tj. zamítáme pro *příliš malé* hodnoty testové statistiky $S_n(\mathbf{X})$;
- $(-\infty, c_L(\alpha)) \cup \langle c_U(\alpha), \infty \rangle$, tj. zamítáme jak pro *příliš malé* tak pro *příliš velké* hodnoty testové statistiky $S_n(\mathbf{X})$;
- $(-\infty, -c_U(\alpha)) \cup \langle c_U(\alpha), \infty \rangle$, tj. zamítáme pro *příliš velké* hodnoty $|S_n(\mathbf{X})|$.

Konstanty $c_L(\alpha)$ a $c_U(\alpha)$, které určují hranice kritického oboru, nazýváme *kritické hodnoty*.[‡] Tyto hodnoty určujeme tak, aby test měl předepsanou hladinu. Jak uvidíme níže, tak ve všech případech, se kterými se setkáme, se dají kritické hodnoty vyjádřit pomocí kvantilů vhodné distribuční funkce G_0 .

Kritický obor tvaru $C(\alpha) = \langle c_U(\alpha), \infty \rangle$

Uvažujme pro jednoduchost nejdříve *přesný test*. Potom kritickou hodnotu $c_U(\alpha)$ volíme tak, aby

$$\sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \geq c_U(\alpha)] = \alpha.$$

My se budeme setkávat pouze s případy, kde lze (snadno) najít $F_0 \in \mathcal{F}_0$ takové, že

$$\sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \geq c] = P_{F_0}[S_n(\mathbf{X}) \geq c] \quad \forall c \in \mathbb{R}. \quad (4.5)$$

* Angl. *unbiased test* † Jak uvidíme později, pro vyjádření kritického oboru pomocí kvantilů by bylo šikovnější uvažovat spíše otevřené kritické obory. To by však trochu pokazilo pozdější dualitu s intervalovými odhady, které se naopak zpravidla uvažují jako otevřené intervaly. ‡ Angl. *critical values*

Toto rozdělení F_0 zpravidla hledáme jako rozdělení, které sice splňuje nulovou hypotézu (tj. je v \mathcal{F}_0), ale je „nejbližší“ alternativě (tj. je „nejblíže“ množině \mathcal{F}_1 , viz také **příklad** (A2) níže). Nechť G_0 nyní značí distribuční funkci $S_n(\mathbf{X})$, pokud rozdělení X_i je F_0 . Potom v případě spojitě distribuční funkce dostáváme $c_U(\alpha) = G_0^{-1}(1 - \alpha)$. V opačném případě je nejjednodušší uvažovat otevřený kritický obor tvaru

$$C(\alpha) = (G_0^{-1}(1 - \alpha), \infty).$$

Všimněme si, že výše uvedený otevřený kritický obor umožňuje správně kontrolovat hladinu, ať už je G_0 spojitá, nebo ne.

Poznámka. V případě, že trváme na uzavřeném kritickém oboru $C(\alpha) = \langle c_U(\alpha), \infty \rangle$ a chceme vyjádřit $c_U(\alpha)$ jako kvantil G_0 , tak se obecný zápis kritické hodnoty $c_U(\alpha)$ poněkud komplikuje. V případě, že G_0 je distribuční funkce diskrétní náhodné veličiny (tj. je po částech konstantní), pak lze zpravidla* volit $c_U(\alpha)$ jako nejbližší vyšší kvantil než $G_0^{-1}(1 - \alpha)$, tj.

$$c_U(\alpha) = \inf \{G_0^{-1}(u) : \text{kde } G_0^{-1}(u) > G_0^{-1}(1 - \alpha) \text{ a } u > 1 - \alpha\}. \quad (4.6)$$

V případě *asymptotického testu* vlastně použijeme jako G_0 distribuční funkci asymptotického rozdělení testové statistiky za nulové hypotézy. Přesněji G_0 je taková distribuční funkce, která splňuje

$$\sup_{F \in \mathcal{F}_0} \lim_{n \rightarrow \infty} P_F[S_n(\mathbf{X}) \geq c] = 1 - G_0(c-) \quad \forall c \in \mathbb{R}.$$

Jelikož v našem případě bude funkce G_0 vždy spojitá, tak pravá strana poslední rovnosti bude tvaru $1 - G_0(c)$ a hledané $c_U(\alpha)$ bude jednoduše $G_0^{-1}(1 - \alpha)$.

Kritický obor tvaru $C(\alpha) = (-\infty, c_L(\alpha))$

Podobně jako výše nechť $F_0 \in \mathcal{F}_0$ je takové rozdělení X_i , že

$$\sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \leq c] = P_{F_0}[S_n(\mathbf{X}) \leq c] \quad \forall c \in \mathbb{R}.$$

Dále nechť G_0 nyní značí distribuční funkci $S_n(\mathbf{X})$, pokud rozdělení X_i je F_0 . Potom v případě spojitě G_0 volíme $c_L(\alpha)$ jako

$$c_L(\alpha) = G_0^{-1}(\alpha).$$

V případě, že G_0 není spojitá, tak je opět nejjednodušší vzít otevřený kritický obor

$$C(\alpha) = (-\infty, c_L(\alpha)) = (-\infty, G_0^{-1}(\alpha)).$$

Tento kritický obor opět správně kontroluje hladinu, ať už G_0 spojitá nebo ne.

Kritický obor tvaru $C(\alpha) = (-\infty, c_L(\alpha)) \cup \langle c_U(\alpha), \infty \rangle$

* Tj. za předpokladu, že níže neděláme infimum z prázdné množiny.

V tomto případě je standardní volit kritické hodnoty $c_L(\alpha)$ a $c_U(\alpha)$ tak, aby

$$\sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \leq c_L(\alpha)] = \sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \geq c_U(\alpha)] = \frac{\alpha}{2}, \quad (4.7)$$

resp.

$$\sup_{F \in \mathcal{F}_0} \lim_{n \rightarrow \infty} P_F[S_n(\mathbf{X}) \leq c_L(\alpha)] = \sup_{F \in \mathcal{F}_0} \lim_{n \rightarrow \infty} P_F[S_n(\mathbf{X}) \geq c_U(\alpha)] = \frac{\alpha}{2}. \quad (4.8)$$

Navíc v situacích, se kterými se budeme setkávat, bude (přesné či asymptotické) rozdělení statistiky $S_n(\mathbf{X})$ za nulové hypotézy vždy stejné, ať už je skutečné rozdělení F z \mathcal{F}_0 jakékoliv (viz také **příklady** (B) a (C) na stranách 81 a 83). Označme toto rozdělení G_0 . Tedy v rovnicích (4.7) a (4.8) můžeme vynechat suprema a podmínka se nám zjednoduší na

$$G_0(c_L(\alpha)) = 1 - G_0(c_U(\alpha)) = \frac{\alpha}{2}.$$

Kritické hodnoty v tomto případě pro G_0 spojitou jsou

$$c_L(\alpha) = G_0^{-1}(\alpha/2) \quad \text{a} \quad c_U(\alpha) = G_0^{-1}(1 - \alpha/2). \quad (4.9)$$

Obecně je pak podobně jako u jednostranných intervalů nejjednodušší uvažovat otevřený kritický obor ve tvaru

$$C(\alpha) = (-\infty, G_0^{-1}(\alpha/2)) \cup (G_0^{-1}(1 - \alpha/2), \infty)$$

Příklad (A2). JEDNOSTRANNÝ TEST STŘEDNÍ HODNOTY NORMÁLNÍHO ROZDĚLENÍ SE ZNÁMÝM ROZPTYLEM.

Máme náhodný výběr X_1, \dots, X_n z rozdělení $F_X = N(\theta_X, \sigma_0^2) \in \mathcal{F}^A = \{N(\theta, \sigma_0^2), \theta \in \mathbb{R}\}$. Testujeme

$$H_0 : \theta_X \leq \theta_0 \quad \text{proti} \quad H_1 : \theta_X > \theta_0.$$

Testová statistika je stejná jako v příkladě A1

$$U_n = \frac{\sqrt{n}(\bar{X}_n - \theta_0)}{\sigma_0}.$$

Její rozdělení pro $\theta_X = \theta_0$ je $N(0, 1)$. Pro hodnoty $\theta_X = \theta_0 + \delta$ máme $U_n \sim N(v_n, 1)$, kde $v_n = \frac{\sqrt{n}\delta}{\sigma_0}$. Je-li porušena hypotéza, pak se rozdělení testové statistiky posouvá do kladných hodnot, a to tím dále, čím větší je n a δ . Příliš velké kladné hodnoty testové statistiky, tedy povedou k zamítnutí hypotézy.

Kritický obor bude mít tvar $C(\alpha) = \langle c_U(\alpha), \infty \rangle$. Kritickou hodnotu $c_U(\alpha)$ určíme tak, aby

$$\sup_{F \in \mathcal{F}_0} P_F[U_n \in C(\alpha)] = \sup_{\theta \in \Theta_0} P_\theta[U_n \in C(\alpha)] = \alpha.$$

Jelikož

$$P_\theta[U_n \in \langle c_U(\alpha), \infty \rangle] = P_\theta\left[\frac{\sqrt{n}(\bar{X}_n - \theta_0)}{\sigma_0} \geq c_U(\alpha)\right] = 1 - \Phi\left(c_U(\alpha) - \frac{\sqrt{n}(\theta - \theta_0)}{\sigma_0}\right)$$

je rostoucí funkce parametru θ , tak

$$\begin{aligned} \sup_{F \in \mathcal{F}_0} P_F[U_n \in C(\alpha)] &= \sup_{\theta: \theta \leq \theta_0} P_\theta[U_n \in \langle c_U(\alpha), \infty \rangle] \\ &= 1 - \Phi\left(c_U(\alpha) - \frac{\sqrt{n}(\theta - \theta_0)}{\sigma_0}\right) = 1 - \Phi(c_U(\alpha)). \end{aligned}$$

Tedy pro $c_U(\alpha) = u_{1-\alpha}$ splňuje tento test podmínku $\sup_{\theta \in \Theta_0} P_\theta[U_n \in C(\alpha)] = \alpha$ a tudíž má hladinu α . Za povšimnutí stojí, že roli F_0 z (4.5), zde hraje rozdělení $N(\theta_0, \sigma_0^2)$. Roli G_0 , tj. distribuční funkce testové statistiky U_n pro X_i s rozdělením F_0 , pak hraje distribuční funkce Φ rozdělení $N(0, 1)$.

Dohromady dostáváme pravidlo

$$\text{zamítání } H_0 : \theta_X \leq \theta_0 \iff U_n = \frac{\sqrt{n}(\bar{X}_n - \theta_0)}{\sigma_0} \geq u_{1-\alpha},$$

tj. zamítáme hypotézu, pokud \bar{X}_n je o více než $\frac{u_{1-\alpha}\sigma_0}{\sqrt{n}}$ větší než θ_0 . Za kvantil $u_{1-\alpha}$ dosazujeme 1,645 pro $\alpha = 0,05$ a 1,282 pro $\alpha = 0,1$. Kritická hodnota pro jednostranný test na hladině α je stejná jako kritická hodnota pro oboustranný test na hladině $\alpha/2$. To je dáno tím, že nyní zamítáme hypotézu pouze v jednom chvostu rozdělení U_n .

Výpočet *silofunkce* je jednodušší než předtím. Vezměme nějaké θ takové, že $\theta - \theta_0 = \delta$ a dostaneme

$$\beta_n(\theta) = P_\theta[U_n \geq u_{1-\alpha}] = P_\theta[U_n - v_n \geq u_{1-\alpha} - v_n] = 1 - \Phi(u_{1-\alpha} - v_n).$$

Průběh silofunkce tohoto testu je zakreslen na obrázku 4.3. Počet pozorování, který je potřeba, aby test dosáhl síly alespoň β proti alternativě $\theta_0 + \delta$, $\delta > 0$, je

$$n \geq (u_{1-\alpha} + u_\beta)^2 \frac{\sigma_0^2}{\delta^2}.$$

Příklad (B). OBOUSTRANNÝ TEST STŘEDNÍ HODNOTY NORMÁLNÍHO ROZDĚLENÍ S NEZNÁMÝM ROZPTYLEM.

Máme náhodný výběr X_1, \dots, X_n z rozdělení $F_X = N(\theta_X, \sigma_X^2) \in \mathcal{F}^B = \{N(\theta, \sigma^2), \theta \in \mathbb{R}, \sigma^2 > 0\}$. Testujeme $H_0 : \theta_X = \theta_0$ proti $H_1 : \theta_X \neq \theta_0$.

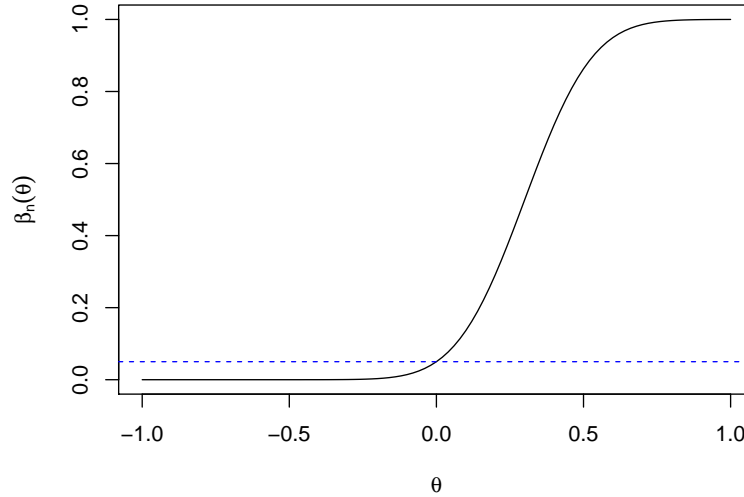
Nemůžeme použít testovou statistiku z příkladů (A1) a (A2), protože neznáme skutečný rozptyl σ_X^2 . Pokud jej však nahradíme výběrovým rozptylem S_n^2 dostaneme statistiku

$$T_n = \frac{\sqrt{n}(\bar{X}_n - \theta_0)}{S_n},$$

kteřá má v tomto modelu za platnosti hypotézy H_0 rozdělení t_{n-1} (viz věta 2.10). Jestliže hypotéza neplatí, tj. $\theta_X - \theta_0 = \delta \neq 0$, pak lze hodnotu této statistiky vyjádřit jako

$$T_n = \frac{Z}{\sqrt{U/(n-1)}},$$

Obrázek 4.3.: Silofunkce testu střední hodnoty normálního rozdělení se známým rozptylem proti pravostranné alternativě pro $\theta_0 = 0$, $\sigma_0^2 = 1$, $n = 30$ a $\alpha = 0,05$.



kde $Z \sim N(v_n, 1)$, $v_n = \frac{\sqrt{n}\delta}{\sigma_X}$, $U \sim \chi_{n-1}^2$ a U , Z jsou nezávislé. Rozdělení této náhodné veličiny se nazývá *necentrální t-rozdělení s $n - 1$ stupni volnosti a parametrem necentrality v_n **. Jeho charakteristiky (hustota, distribuční funkce, momenty) mají komplikovaný tvar, ale stačí vědět, že pro velké n jej lze aproximovat rozdělením $N(v_n, 1)$.

I zde tedy platí, že je-li porušena hypotéza, pak se rozdělení testové statistiky posouvá pryč od nuly, a to tím dále, čím větší je n a $|\theta_X - \theta_0|$. Hodnoty testové statistiky, které jsou daleko od nuly, tedy povedou k zamítnutí hypotézy.

Kritický obor bude mít tvar $(-\infty, c_L(\alpha)) \cup (c_U(\alpha), \infty)$. Všimněme si, že roli G_0 zde hraje rozdělení t_{n-1} , neboť za H_0 platí $T_n \sim t_{n-1}$, ať už má σ^2 jakékoliv (kladnou) hodnotu. Vzhledem k symetrii rozdělení t_{n-1} tak s využitím (4.9) dostáváme

$$c_L(\alpha) = t_{n-1}(\alpha/2) = -t_{n-1}(1 - \alpha/2), \quad c_U(\alpha) = t_{n-1}(1 - \alpha/2).$$

Ze cvičných důvodů ještě ověříme, že test má při takové volbě kritických hodnot opravdu hladinu α . Počítejme

$$\sup_{F \in \mathcal{F}_0} P_F(T_n \in C(\alpha)) = \sup_{\sigma^2 > 0} P_{\theta_0, \sigma^2}(T_n \leq -t_{n-1}(1 - \alpha/2) \text{ nebo } T_n \geq t_{n-1}(1 - \alpha/2)) = \alpha.$$

Test tedy má přesně hladinu α a dostáváme pravidlo

$$\text{zamítni } H_0 : \theta_X = \theta_0 \iff |T_n| = \left| \frac{\sqrt{n}(\bar{X}_n - \theta_0)}{S_n} \right| \geq t_{n-1}(1 - \alpha/2).$$

* Angl. *non-central t distribution with $n - 1$ degrees of freedom and noncentrality parameter v_n*

To znamená, že hypotéza bude zamítnuta, pokud se bude průměr \bar{X}_n lišit od hypotetické hodnoty θ_0 o více než $\frac{t_{n-1}(1-\alpha/2)S_n}{\sqrt{n}}$. Tento test se nazývá *jednovýběrový t-test**.

Silofunkci získáme podobným postupem jako v příkladě (1A). Vezměme nějaké θ takové, že $\theta - \theta_0 = \delta \neq 0$. Pokud θ je skutečný parametr, pak rozdělení T_n je necentrální t s $n - 1$ stupni volnosti a parametrem necentrality $v_n = \frac{\sqrt{n}\delta}{\sigma_X}$. Označme distribuční funkci tohoto rozdělení G_{n,v_n} a počítejme

$$\begin{aligned}\beta_n(\theta, \sigma_X^2) &= P_{\theta, \sigma_X^2}[T_n \in C(\alpha)] \\ &= P_{\theta, \sigma_X^2}[T_n \leq -t_{n-1}(1 - \alpha/2)] + P_{\theta, \sigma_X^2}[T_n \geq t_{n-1}(1 - \alpha/2)] \\ &= G_{n,v_n}(-t_{n-1}(1 - \alpha/2)) + 1 - G_{n,v_n}(t_{n-1}(1 - \alpha/2)).\end{aligned}$$

Necentrální t -rozdělení nemá symetrickou hustotu, takže výsledek již nejde dále upravovat. Pokud je počet pozorování n dostatečně velký, můžeme aproximovat sílu pomocí vzorce (4.2) nebo (4.3).

Ze vzorce (4.3) lze získat aproximaci pro počet pozorování n potřebný k tomu, aby test dosáhl síly alespoň β . Požadovaný rozsah výběru je

$$n \geq (u_{1-\alpha/2} + u_\beta)^2 \frac{\sigma_X^2}{\delta^2} + 1,$$

Jednička se k výsledku přidává proto, aby trochu zkompenzovala nahrazení t -rozdělení normálním. K výpočtu síly a rozsahu výběru je třeba znát skutečný rozptyl σ_X^2 nebo jej nahradit nějakým předběžným odhadem (tyto výpočty obvykle provádíme předtím, než získáme data).

Příklad (C). OBOUSTRANNÝ TEST STŘEDNÍ HODNOTY LIBOVOLNÉHO ROZDĚLENÍ S KONEČNÝM ROZPTYLEM.

Máme náhodný výběr X_1, \dots, X_n z rozdělení $F_X \in \mathcal{F}^C = \mathcal{L}_+^2$. Označme $E X_i = \theta_X$, $\text{var } X_i = \sigma_X^2$. Testujeme $H_0 : \theta_X = \theta_0$ proti $H_1 : \theta_X \neq \theta_0$.

Podle věty 2.9 (limitní věta o T_n) má v tomto modelu náhodná veličina

$$T_n = \frac{\sqrt{n}(\bar{X}_n - \theta_0)}{S_n},$$

za platnosti hypotézy H_0 asymptoticky rozdělení $N(0, 1)$. Jestliže hypotéza neplatí, tj. $\theta_X - \theta_0 = \delta \neq 0$, pak se dá snadno ukázat[†], že testová statistika

$$T_n = \frac{\sqrt{n}(\bar{X}_n - \theta_X + \theta_X - \theta_0)}{S_n} = \frac{\sqrt{n}(\bar{X}_n - \theta_X)}{S_n} + \sqrt{n} \frac{\delta}{S_n}$$

konverguje v pravděpodobnosti do $+\infty$ nebo $-\infty$ podle toho, jaké znaménko má δ . Hodnoty testové statistiky, které jsou daleko od nuly, tedy povedou k zamítnutí hypotézy.

* Angl. *one-sample t-test* † Doporučuji si rozmyslet jako cvičení.

Kritický obor bude mít tvar $(-\infty, c_L(\alpha)) \cup (c_U(\alpha), \infty)$. Všimněme si, že

$$\sup_{F \in \mathcal{F}_0} \lim_{n \rightarrow \infty} P_F(|T_n| \geq u_{1-\alpha/2}) = P(|Z| \geq u_{1-\alpha/2}) = \alpha,$$

kde $Z \sim N(0, 1)$. Tedy roli G_0 v (4.9) hraje distribuční funkce rozdělení $N(0, 1)$. Tudiž kritické hodnoty $c_U(\alpha) = -c_L(\alpha) = u_{1-\alpha/2}$ zaručují, že hladina testu je asymptoticky rovna α .

Místo kritické hodnoty $u_{1-\alpha/2}$ můžeme použít $t_{n-1}(1 - \alpha/2)$, protože provádíme asymptotický test a $t_{n-1}(1 - \alpha/2) \rightarrow u_{1-\alpha/2}$ pro $n \rightarrow \infty$. Jelikož $|t_{n-1}(\alpha)| \geq |u_\alpha|$, tak test bude s využitím kvantilů t -rozdělení konservativnější, než kdybychom použili kvantily normovaného normálního rozdělení.

Celkem tedy dostáváme pravidlo

$$\text{zamítni } H_0 : \theta_X = \theta_0 \iff |T_n| = \left| \frac{\sqrt{n}(\bar{X}_n - \theta_0)}{S_n} \right| \geq t_{n-1}(1 - \alpha/2).$$

Jedná se tedy opět o jednovýběrový t -test. Ukázali jsme, že jakožto asymptotický test jej můžeme použít pro libovolná data s konečným rozptylem.

Cvičení.

1. V příkladu A1 (str. 75) uvažujte test $(U_n, C(\alpha))$, kde $C(\alpha) = \langle u_{1/2-\alpha/2}, u_{1/2+\alpha/2} \rangle$. Ukažte, že tento test má hladinu přesně α . Dále ukažte, že pro tento test platí, že pro všechna θ různá od θ_0

$$\beta_n(\theta) < \alpha \quad \text{a} \quad \lim_{n \rightarrow \infty} \beta_n(\theta) = 0.$$

2. V příkladu A1 (str. 75) uvažujte test $(U_n, C(\alpha))$, kde $C(\alpha) = \langle u_{1-\alpha}, \infty \rangle$. Ukažte, že tento test má hladinu přesně α . Je tento test nestranný? Pro jaká θ je tento test konzistentní?
3. Dokažte o testu z příkladu A2 (str. 80), že je nestranný a konzistentní.
4. Dokažte o testu z příkladu B (str. 81), že je nestranný a konzistentní.
Návod. Pro důkaz nestrannosti můžete využít toho, že pro náhodnou veličinu Z s necentrálním t -rozdělením se stupni volnosti v a nenulovým parametrem necentrality platí $P(|Z_n| \geq t_v(1 - \alpha/2)) > \alpha$.
5. Dokažte o testu z příkladu C (viz předchozí strana), že je konzistentní.
6. Oddělení PR jedné střední školy by rádo prokázalo, že střední hodnota IQ jejich studentů je vyšší než 105. Přičemž očekávají, že skutečná střední hodnota IQ jejich studentů je 110 a směrodatná chyba rozdělení IQ těchto studentů je 15. Spočítejte, kolika studentům je třeba změřit IQ, aby za těchto předpokladů test na hladině 5% s pravděpodobností 95% prokázal, že střední hodnota IQ studentů dané školy je vyšší než 105.

4.3. P-HODNOTA

Posuzovat výsledek testu podle toho, zda $S_n(\mathbf{X})$ padne do C , není jediný ani nejběžnější způsob vyhodnocování testů. Výsledek testu se v praxi nejčastěji posuzuje pomocí tzv. **p-hodnoty**. Ta odpovídá *nejmenší možné hladině, na které bychom ještě zamítli nulovou hypotézu*.

Uvažujme tedy hypotézu $H_0 : \theta_X \in \Theta_0$ proti alternativě $H_1 : \theta_X \in \Theta_1$ a nechť pro dané $\alpha \in (0, 1)$ je $(S_n(\mathbf{X}), C(\alpha))$ test s předepsanou hladinou α . Pro určitost ještě dodefinujeme $C(1) = (-\infty, \infty)$.

Jak víme z poznámky na straně 73, v případě že $S_n(\mathbf{X})$ má za nulové hypotézy diskrétní rozdělení, pak není možné dosáhnout zcela libovolné hladiny testu. Proto pokud daná hladina α je nedosažitelná, tak budeme jako $C(\alpha)$ značit kritický obor testu, který dává hladinu $\alpha' < \alpha$, kde α' je nejbližší k původně požadovanému α .

Definice 4.8 (P-hodnota) Nechť $s_x = S_n(x)$ je realizovaná hodnota testové statistiky. Pak *p-hodnotu** neboli *dosaženou hladinu testu* definujeme jako

$$p(x) = \inf \{ \alpha \in (0, 1) : s_x \in C(\alpha) \}.$$

Je-li test $(S_n(\mathbf{X}), C(\alpha))$ přesný (resp. asymptotický), tak p-hodnotu nazýváme *přesnou* (resp. *asymptotickou*).

Má-li test předepsanou hladinu α , tak se můžeme jednoduše rozhodovat podle následujícího pravidla

$$\begin{aligned} H_0 \text{ zamítáme, jestliže } p(x) \leq \alpha, \\ H_0 \text{ nezamítáme, jestliže } p(x) > \alpha. \end{aligned} \tag{4.10}$$

Tudíž známe-li p-hodnotu $p(x)$, pak můžeme hypotézu zamítnout na všech hladinách $\alpha' \geq p(x)$, ale nemůžeme ji zamítnout na hladinách $\alpha' < p(x)$. Proto se p-hodnotě říká *dosažená hladina testu*.

Zamítáme-li pomocí p-hodnoty, **nemusíme uvádět kritický obor** a nemusíme jej přepočítávat, pokud se rozhodneme změnit hladinu testu. Upozorníme však, že *měnit hladinu testu poté, co je znám výsledek, není legitimní*.

Poznámka.

- P-hodnotu můžeme chápat jako *míru souladu dat* s hypotézou. Pokud $p(x) \ll \alpha$, data zamítají hypotézu s velkou „rezervou“. Pokud je $p(x)$ blízké α , tak se zase někdy říká, že výsledek je „na hraně statistické významnosti“.
- P-hodnotu **není možné** vykládat jako „pravděpodobnost, že nulová hypotéza platí“. Platnost nulové hypotézy totiž není náhodný, ale deterministický jev.
- Pokud má $S_n(\mathbf{X})$ *diskrétní* rozdělení, pak pravidlo (4.10) dává test, který má nejbližší možnou dosažitelnou hladinu α' takovou, že $\alpha' \leq \alpha$.

* Angl. *p-value*

4.3.1. VÝPOČET P-HODNOTY PRO JEDNOSTRANNÝ KRITICKÝ OBOR

Jak je vidět z definice 4.8, tak p-hodnota je funkcí realizovaných dat x a její výpočet závisí na použité testové statistice $S_n(\mathbf{X})$ a způsobu, jak se mění kritický obor $C(\alpha)$, když měníme α . Nejjednodušším případem je situace, kdy kritický obor je jednostranný, tj. zamítáme pro příliš velké (malé) hodnoty testové statistiky.

Předpokládejme nejdříve, že zamítáme pro příliš velké hodnoty testové statistiky, tj. $C(\alpha) = \langle c_U(\alpha), \infty \rangle$. Z kapitoly 4.2.3 víme, že $c_U(\alpha) = G_0^{-1}(1 - \alpha)$, přičemž obecně je nejjednodušší uvažovat otevřený kritický obor $C(\alpha) = (G_0^{-1}(1 - \alpha), \infty)$. Připomeňme, že G_0 je v tomto případě distribuční funkce taková, že

$$\sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \geq c] = 1 - G_0(c-) \quad \forall c \in \mathbb{R}.$$

V tomto případě z definice p-hodnoty tudíž dostáváme

$$p(x) = \inf \{ \alpha \in (0, 1) : s_x > G_0^{-1}(1 - \alpha) \} = 1 - G_0(s_x -). \quad (4.11)$$

Naprosto analogicky lze postupovat v případě, že zamítáme pro příliš malé hodnoty testové statistiky, tj. v případě kritického oboru $C(\alpha) = (-\infty, c_L(\alpha))$. Nechť G_0 je distribuční funkce taková, že

$$\sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \leq c] = G_0(c) \quad \forall c \in \mathbb{R}.$$

Potom dostáváme, že $c_L(\alpha) = G_0^{-1}(\alpha)$, přičemž pro nespojitou G_0 je nejjednodušší uvažovat otevřený kritický obor $C(\alpha) = (-\infty, G_0^{-1}(\alpha))$. Tedy z definice p-hodnoty

$$p(x) = \inf \{ \alpha \in (0, 1) : s_x < G_0^{-1}(\alpha) \} = G_0(s_x). \quad (4.12)$$

Poznámka.

- Vzorce (4.11) a (4.12) lze použít i pro *asymptotickou p-hodnotu* v případě, že G_0 je distribuční funkce asymptotického rozdělení testové statistiky za hypotézy. Tj. uvažujme například kritický obor tvaru $C(\alpha) = \langle c_U(\alpha), \infty \rangle$. Pak potřebujeme, aby G_0 byla taková distribuční funkce, která splňuje

$$\sup_{F \in \mathcal{F}_0} \lim_{n \rightarrow \infty} P_F[S_n(\mathbf{X}) \geq c] = 1 - G_0(c-), \quad \forall c \in \mathbb{R}.$$

Podobně pro kritický obor $C(\alpha) = (-\infty, c_L(\alpha))$ potřebujeme, aby G_0 splňovala

$$\sup_{F \in \mathcal{F}_0} \lim_{n \rightarrow \infty} P_F[S_n(\mathbf{X}) \leq c] = G_0(c), \quad \forall c \in \mathbb{R}.$$

- Všimněme si, že pro kritický obor $C(\alpha) = \langle c_U(\alpha), \infty \rangle$ můžeme vzorec (4.11) pro p-hodnotu přepsat do tvaru

$$p(x) = 1 - G_0(s_x -) = \sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \geq s_x].$$

Podobně pro kritický obor tvaru $C(\alpha) = (-\infty, c_L(\alpha))$

$$p(x) = G_0(s_x) = \sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \leq s_x].$$

Na p-hodnotu je tedy také možné pohlížet jako na (maximální možnou) pravděpodobnost, že bychom za platnosti hypotézy napozorovali data, která by byla s hypotézou *ve stejném nebo větším rozporu*, než data, která analyzujeme.

Příklad (A). TEST STŘEDNÍ HODNOTY NORMÁLNÍHO ROZDĚLENÍ SE ZNÁMÝM ROZPTYLEM.

Máme náhodný výběr X_1, \dots, X_n z rozdělení $F_X = N(\theta_X, \sigma_0^2) \in \mathcal{F}^A = \{N(\theta, \sigma_0^2), \theta \in \mathbb{R}\}$. Testujeme

$$H_0 : \theta_X \geq \theta_0 \quad \text{proti} \quad H_1 : \theta_X < \theta_0.$$

Testová statistika je v tomto případě

$$U_n = \frac{\sqrt{n}(\bar{X}_n - \theta_0)}{\sigma_0},$$

přičemž zamítáme pro *malé* hodnoty testové statistiky.

Všimněme si (viz příklad 4.2.2), že rozdělení testové statistiky je $N(v_n, 1)$, kde střední hodnota $v_n = \frac{\sqrt{n}(\theta_X - \theta_0)}{\sigma_0}$ je za nulové hypotézy nezáporná. Nechť u_x je napozorovaná hodnota testové statistiky U_n . Jelikož kritický obor je tvaru $C(\alpha) = (-\infty, c_L(\alpha))$, kde $c_L(\alpha) = \Phi^{-1}(\alpha)$, tak z definice 4.8 dostáváme

$$p(x) = \inf\{\alpha \in (0, 1) : u_x \leq \Phi^{-1}(\alpha)\} = \Phi(u_x).$$

Příklad (B). TEST STŘEDNÍ HODNOTY NORMÁLNÍHO ROZDĚLENÍ S NEZNÁMÝM ROZPTYLEM.

Máme náhodný výběr X_1, \dots, X_n , $n = 26$, z rozdělení $F_X \in \mathcal{F}^B$ a uvažujme $\theta_X = E X_i$. Testujeme

$$H_0 : \theta_X \leq \theta_0 \quad \text{proti} \quad H_1 : \theta_X > \theta_0.$$

K tomu použijeme statistiku

$$T_n = \frac{\sqrt{n}(\bar{X}_n - \theta_0)}{S_n},$$

přičemž nulovou hypotézu zamítáme, pokud je hodnota testové statistiky příliš *velká*.

Předpokládejme, že jsme spočítali testovou statistiku a její výsledek označme jako t_x . V příkladu B na straně 81 bylo ukázáno, že testová statistika T_n má necentrální t -rozdělení s $n - 1$ stupni volnosti a parametrem necentrality $v_n = \frac{\sqrt{n}(\theta_X - \theta_0)}{\sigma_X}$. Všimněme si, že za nulové hypotézy je tento parametr záporný nebo nulový.

Proti nulové hypotéze budou tedy svědčit velké hodnoty testové statistiky a kritický obor bude tvaru $C(\alpha) = (c_U(\alpha), \infty)$. Dále vzhledem k výše uvedenému bude roli G_0 při

určování kritického oboru pomocí vzorce (4.6) hrát distribuční funkce rozdělení t_{n-1} , tj. v našem případě t_{25} . Tedy $c_U(\alpha) = t_{25}(1 - \alpha)$ a tudíž

$$p(x) = \inf\{\alpha \in (0, 1) : t_x \geq t_{25}(1 - \alpha)\} = 1 - G_{25}(t_x),$$

kde G_{25} je distribuční funkce rozdělení t_{25} .

4.3.2. VÝPOČET P-HODNOTY PRO OBOUSTRANNÝ KRITICKÝ OBOR

Z kapitoly 4.2.3 víme, že v následujícím se budeme setkávat pouze se situacemi, kdy přesné (či asymptotické) rozdělení testové statistiky $S_n(\mathbf{X})$ nezávisí na volbě F z \mathcal{F}_0 . Označme distribuční funkci tohoto (přesného či asymptotického) rozdělení $S_n(\mathbf{X})$ jako G_0 . Potom dle (4.9) plyne, že

$$c_L(\alpha) = G_0^{-1}(\alpha/2), \quad \text{a} \quad c_U(\alpha) = G_0^{-1}(1 - \alpha/2),$$

přičemž pro jednoduchost zápisu budeme v následujícím uvažovat otevřený kritický obor

$$C(\alpha) = (-\infty, G_0^{-1}(\alpha/2)) \cup (G_0^{-1}(1 - \alpha/2), \infty).$$

Tudíž z definice p-hodnoty

$$\begin{aligned} p(x) &= \inf\{\alpha \in (0, 1) : s_x < G_0^{-1}(\alpha/2) \text{ nebo } s_x > G_0^{-1}(1 - \alpha/2)\} \\ &= 2 \min\{1 - G_0(s_x-), G_0(s_x)\}. \end{aligned} \quad (4.13)$$

Všimněme si, že

$$p(x) = 2 \min\{1 - G_0(s_x-), G_0(s_x)\} = 2 \min\{P_{F_0}[S_n(\mathbf{X}) \leq s_x], P_{F_0}[S_n(\mathbf{X}) \geq s_x]\},$$

kde poslední vzorec může být snadnější na zapamatování i interpretaci o souladu dat s nulovou hypotézou.

Vzorec (4.13) můžeme dále zjednodušit v případě, že přesné (resp. asymptotické) rozdělení G_0 je symetrické kolem 0, tj.

$$1 - G_0(|s_x| -) = G_0(|s_x|).$$

Potom

$$p(x) = 2(1 - G_0(|s_x| -)). \quad (4.14)$$

Za povšimnutí také stojí, že v tomto případě se dá p-hodnota vyjádřit jako

$$\begin{aligned} p(x) &= 2 \min\{1 - G_0(s_x-), G_0(s_x)\} \\ &= 2 \min\{P_{F_0}[S_n(\mathbf{X}) \leq s_x], P_{F_0}[S_n(\mathbf{X}) \geq s_x]\} \\ &= 2 P_{F_0}[S_n(\mathbf{X}) \geq |s_x|] = P_{F_0}[|S_n(\mathbf{X})| \geq |s_x|]. \end{aligned}$$

Tj. podobně jako u jednostranných kritických oborů, tak i zde můžeme na p-hodnotu pohlížet jako na pravděpodobnost, že bychom za platnosti hypotézy napozorovali data, která by byla s hypotézou *ve stejném nebo větším rozporu*, než data, která analyzujeme.

Příklad (A). TEST STŘEDNÍ HODNOTY NORMÁLNÍHO ROZDĚLENÍ SE ZNÁMÝM ROZPTY-
LEM.

Máme náhodný výběr X_1, \dots, X_n z rozdělení $F_X = N(\theta_X, \sigma_0^2) \in \mathcal{F}^A = \{N(\theta, \sigma_0^2), \theta \in \mathbb{R}\}$ a zajímá nás test hypotézy

$$H_0 : \theta_X = \theta_0 \quad \text{proti} \quad H_1 : \theta_X \neq \theta_0,$$

Testová statistika je v tomto případě

$$U_n = \frac{\sqrt{n} (\bar{X}_n - \theta_0)}{\sigma_0},$$

a zamítáme pro příliš velké, resp. příliš malé hodnoty testové statistiky. Výpočet p-hodnoty je v tomto případě jednoduchý, protože hypotéza obsahuje právě jedno rozdělení $N(\theta_0, \sigma_0^2)$, které bude hrát roli rozdělení F_0 . Navíc v tomto případě má za nulové hypotézy testová statistika U_n rozdělení $N(0, 1)$, které je symetrické kolem nuly. Tudíž p-hodnota testu se spočte jako

$$p(\mathbf{x}) = 2 \min \{1 - \Phi(u_{\mathbf{x}}), \Phi(u_{\mathbf{x}})\} = 2 (1 - \Phi(|u_{\mathbf{x}}|)).$$

Příklad (B). TEST STŘEDNÍ HODNOTY NORMÁLNÍHO ROZDĚLENÍ S NEZNÁMÝM ROZPTY-
LEM.

Máme náhodný výběr X_1, \dots, X_n , $n = 26$, z rozdělení $F_X \in \mathcal{F}^B$ a uvažujme $\theta_X = E X_i$. Testujeme $H_0 : \theta_X = \theta_0$ proti $H_1 : \theta_X \neq \theta_0$. K tomu použijeme statistiku

$$T_n = \frac{\sqrt{n} (\bar{X}_n - \theta_0)}{S_n},$$

přičemž nulovou hypotézu zamítáme, pokud je hodnota testové statistiky příliš malá nebo příliš velká.

Předpokládejme, že jsme spočítali testovou statistiku a její výsledek označme jako $t_{\mathbf{x}}$. V příkladu B na straně 81 bylo ukázáno, že testová statistika T_n má necentrální t -rozdělení s $n - 1$ stupni volnosti a parametrem necentrality $v_n = \sqrt{n}(\theta_X - \theta_0)/\sigma_X$. Všimněme si, že za nulové hypotézy je tento parametr nulový. Tedy

$$p(\mathbf{x}) = 2 \min \{1 - G_{25}(t_{\mathbf{x}}), G_{25}(t_{\mathbf{x}})\} = 2 (1 - G_{25}(|t_{\mathbf{x}}|)),$$

kde jsme využili toho, že t -rozdělení s $n - 1$ stupni volnosti je symetrické kolem nuly.

Speciálně pro $t_{\mathbf{x}} = 1,37$ pak dostáváme

$$p(\mathbf{x}) = 2 (1 - G_{25}(|1,37|)) \doteq 0,183.$$

Příklad (C). Máme náhodný výběr X_1, \dots, X_n , $n = 26$, z rozdělení $F_X \in \mathcal{F}^C = \mathcal{L}_+^2$ se střední hodnotou $E X_i = \theta_X$. Testujeme $H_0 : \theta_X = \theta_0$ proti $H_1 : \theta_X \neq \theta_0$. Testová statistika T_n má za platnosti hypotézy přibližně rozdělení $N(0, 1)$, které je symetrické

kolem 0. Spočítali jsme testovou statistiku a její výsledek je $t_x = 1,37$. Asymptotická p-hodnota pro tento test se spočítá pomocí (4.14) jako

$$p(x) = 2(1 - \Phi(|1,37|)) \doteq 0,171. \quad (4.15)$$

Testujeme-li na hladině $\alpha = 0,05$, nemůžeme zamítnout hypotézu, neboť $p(x) > 0,05$. Kdybychom si však před provedením testu stanovili hladinu $\alpha' = 0,2$, hypotézu bychom zamítnout mohli.

Všimněme si, že v modelu \mathcal{F}^B (tj. množina normálních rozdělení s neznámých rozptylem) by přesná p-hodnota s využitím (4.14) byla

$$p(x) = 2(1 - G_{25}(|1,37|)) \doteq 0,183, \quad (4.16)$$

kde G_{25} značí distribuční funkci t -rozdělení t_{25} . Jelikož tato p-hodnota je vyšší než asymptotická p-hodnota (4.15), tak se i v modelu \mathcal{F}^C z důvodů opatrnosti (konzervativnosti) používá p-hodnota (4.16) spočtená pomocí rozdělení t_{25} . Jelikož rozdělení t_{n-1} konverguje (v distribuci) k normálnímu rozdělení $N(0, 1)$, tak můžeme i na p-hodnotu (4.16) nahlížet jako na asymptotickou p-hodnotu pro model \mathcal{F}^C .

Za povšimnutí stojí, že vzorec (4.16) dostaneme také přímo z definice p-hodnoty, když jako kritické hodnoty použijeme $c_L(\alpha) = t_{n-1}(\alpha/2)$ a $c_U(\alpha) = t_{n-1}(1 - \alpha/2)$. V tomto případě totiž

$$\begin{aligned} p(x) &= \inf \{ \alpha \in (0, 1) : 1,37 \leq t_{n-1}(\alpha/2) \text{ nebo } 1,37 \geq t_{n-1}(1 - \alpha/2) \} \\ &= 2 \min \{ 1 - G_{25}(1,37), G_0(1,37) \} = 2 [1 - G_{25}(|1,37|)]. \end{aligned}$$

4.3.3. ROZDĚLENÍ P-HODNOTY ZA NULOVÉ HYPOTÉZY

Uvažujme nyní p-hodnotu $p(\mathbf{X})$ jakožto náhodnou veličinu, čili statistiku spočítanou z náhodného výběru \mathbf{X} . Lze ukázat, že za určitých předpokladů má $p(\mathbf{X})$ za platnosti nulové hypotézy rovnoměrné rozdělení na intervalu $(0, 1)$.

Tvrzení 4.1 Nechť platí hypotéza (tj. $F_X \in \mathcal{F}_0$). Nechť navíc pro (skutečné rozdělení) F_X platí, že

$$\sup_{F \in \mathcal{F}_0} P_F[S_n(\mathbf{X}) \in C(\alpha)] = P_{F_X}[S_n(\mathbf{X}) \in C(\alpha)], \forall \alpha \in (0, 1). \quad (4.17)$$

Dále nechť má testová statistika $S_n(\mathbf{X})$ **spojité** rozdělení. Pak $p(\mathbf{X}) \sim R(0, 1)$.

Důkaz. Označme $U = G_0(S_n(\mathbf{X}))$, kde G_0 je distribuční funkce náhodné veličiny $S_n(\mathbf{X})$, pokud rozdělení X_i je F_X . Všimněme si, že v tomto případě má náhodná veličina U rovnoměrné rozdělení na $(0, 1)$ (viz lemma A.5). Dále nechť s_x značí napozorovanou hodnotu statistiky $S_n(\mathbf{X})$.

Tvrzení dále dokážeme podle jednotlivých druhů kritických oborů.

(i) $C(\alpha) = \langle c_U(\alpha), \infty \rangle$

V tomto případě s využitím (4.11) pro p-hodnotu platí $p(x) = 1 - G_0(s_x)$. A tedy pro distribuční funkci náhodné veličiny $p(\mathbf{X})$ platí

$$P_{F_X}[p(\mathbf{X}) \leq u] = P_{F_X}[1 - G_0(S_n(\mathbf{X})) \leq u] = P[1 - U \leq u] = P[1 - u \leq U] = u,$$

pro $\forall u \in (0, 1)$. Tedy distribuční funkce $p(\mathbf{X})$ se shoduje s distribuční funkcí rovnoměrného rozdělení na $(0, 1)$, což bylo dokázat.

(ii) $C(\alpha) = (-\infty, c_L(\alpha))$

V tomto případě s využitím (4.12) pro $\forall u \in (0, 1)$

$$P_{F_X}[p(\mathbf{X}) \leq u] = P_{F_X}[G_0(S_n(\mathbf{X})) \leq u] = P[U \leq u] = u.$$

(iii) $C(\alpha) = (-\infty, c_L(\alpha)) \cup (c_U(\alpha), \infty)$

S využitím (4.13) pro $\forall u \in (0, 1)$

$$\begin{aligned} P_{F_X}[p(\mathbf{X}) \leq u] &= P_{F_X}[2 \min\{1 - G_0(S_n(\mathbf{X})), G_0(S_n(\mathbf{X}))\} \leq u] \\ &= P[2 \min\{1 - U, U\} \leq u] \\ &= P[2 \min\{1 - U, U\} \leq u, U \leq \frac{1}{2}] + P[2 \min\{1 - U, U\} \leq u, U \geq \frac{1}{2}] \\ &= P[2U \leq u, U \leq \frac{1}{2}] + P[2(1 - U) \leq u, U \geq \frac{1}{2}] \\ &= P[U \leq \min\{\frac{u}{2}, \frac{1}{2}\}] + P[U \geq \max\{1 - \frac{u}{2}, \frac{1}{2}\}] \\ &= \frac{u}{2} + 1 - (1 - \frac{u}{2}) = u. \end{aligned}$$

□

Poznámka. Předchozí tvrzení neplatí, pokud je rozdělení testové statistiky diskrétní. Tvrzení by také neplatilo, pokud by sice platila hypotéza (tj. $F_X \in \mathcal{F}_0$), ale F_X by nebylo „nejbližší“ alternativě, (tj. v (4.17) bychom nemohli nahradit $\sup_{F \in \mathcal{F}_0} P_F$ za P_{F_X}).

4.4. DUALITA INTERVALOVÝCH ODHADŮ A TESTOVÁNÍ HYPOTÉZ

Uvažujme náhodný výběr $\mathbf{X} = (X_1, \dots, X_n)$ z rozdělení $F_X \in \mathcal{F}$, kde \mathcal{F} je model. Nechť $\theta = t(F) \in \mathbb{R}$ je parametr a $\theta_X = t(F_X)$ je jeho skutečná hodnota. V kapitole 3.4 jsme řešili problém intervalového odhadu parametru θ_X , tj. hledali jsme náhodné veličiny $\eta_L(\mathbf{X})$ a $\eta_U(\mathbf{X})$ takové, že

$$P_{F_X}[(\eta_L(\mathbf{X}), \eta_U(\mathbf{X})) \ni \theta_X] = 1 - \alpha \quad (\text{nebo } \xrightarrow{n \rightarrow \infty} 1 - \alpha) \quad \text{pro } \forall F_X \in \mathcal{F}.$$

V této kapitole se zabýváme testováním hypotéz, speciálně hypotézy

$$H_0 : \theta_X = \theta_0 \text{ proti } H_1 : \theta_X \neq \theta_0.$$

Oba problémy se řeší postupy, které se v určitých rysech shodují, ale liší se v detailech.

Následující věta ukazuje, že mezi problémem testování hypotézy o parametru a problémem hledání intervalového odhadu pro ten samý parametr existuje jakási dualita. Intervalový odhad můžeme použít k testování hypotéz a test hypotézy můžeme převést na intervalový odhad.

Tvrzení 4.2 (Dualita intervalových odhadů a testování)

- (i) Nechť je dán oboustranný interval spolehlivosti pro parametr θ_X s pravděpodobností pokrytí $1-\alpha$ (přesnou nebo asymptotickou), který má tvar $(\eta_L(\mathbf{X}), \eta_U(\mathbf{X}))$. Uvažujme test hypotézy $H_0 : \theta_X = \theta_0$ proti $H_1 : \theta_X \neq \theta_0$ založený na rozhodovacím pravidle

$$\begin{aligned} H_0 \text{ zamítáme, jestliže } \theta_0 \notin (\eta_L(\mathbf{X}), \eta_U(\mathbf{X})) \\ H_0 \text{ nezamítáme, jestliže } \theta_0 \in (\eta_L(\mathbf{X}), \eta_U(\mathbf{X})). \end{aligned} \quad (4.18)$$

Pak tento test má hladinu α (přesně nebo asymptoticky).

- (ii) Nechť pro všechna $\theta \in \Theta$ je dán test $(S_n(\mathbf{X}, \theta), C_\theta(\alpha))$ hypotézy $H_0 : \theta_X = \theta$ proti $H_1 : \theta_X \neq \theta$ takový, že pro všechna F pro která $\theta = t(F)$

$$P_F[S_n(\mathbf{X}, \theta) \in C_\theta(\alpha)] = \alpha \quad (\text{nebo } \xrightarrow{n \rightarrow \infty} \alpha).$$

Sestavme množinu $B_n(\mathbf{X})$ obsahující všechny parametry $\theta \in \Theta$, pro něž se při pozorovaných datech \mathbf{X} nezamítá hypotéza $H_0 : \theta_X = \theta$. Pak pro všechna $F_X \in \mathcal{F}$

$$P_{F_X}[B_n(\mathbf{X}) \ni \theta_X] = 1 - \alpha \quad (\text{nebo } \xrightarrow{n \rightarrow \infty} 1 - \alpha).$$

Je-li $B_n(\mathbf{X})$ interval, pak se jedná o interval spolehlivosti pro parametr θ_X s pravděpodobností pokrytí $1 - \alpha$ (přesnou nebo asymptotickou).

Důkaz. Část (i) Nechť $(\eta_L(\mathbf{X}), \eta_U(\mathbf{X}))$ je přesný interval spolehlivosti. Pro asymptotický interval by byl důkaz analogický.

Interval spolehlivosti pro skutečnou hodnotu parametru θ_X splňuje

$$P_{F_X}[(\eta_L(\mathbf{X}), \eta_U(\mathbf{X})) \ni \theta_X] = 1 - \alpha, \quad \text{pro } \forall F_X \in \mathcal{F}.$$

Tedy za platnosti nulové hypotézy, tj. pro $\theta_X = \theta_0$, pro všechna $F \in \mathcal{F}_0 = \{F \in \mathcal{F} : t(F) = \theta_0\}$ platí

$$P_F[(\eta_L(\mathbf{X}), \eta_U(\mathbf{X})) \ni \theta_0] = 1 - \alpha.$$

Tedy hladina testu daného předpisem (4.18) je

$$\sup_{F \in \mathcal{F}_0} P_F[(\eta_L(\mathbf{X}), \eta_U(\mathbf{X})) \not\ni \theta_0] = \alpha,$$

což bylo dokázat.

Část (ii) Nechť pro všechna $\theta \in \Theta$ je $(S_n(\mathbf{X}, \theta), C_\theta(\alpha))$ přesný test hypotézy $H_0 : \theta_X = \theta$ proti alternativě $H_1 : \theta_X \neq \theta$ s hladinou α . Pro asymptotický test by byl důkaz obdobný.

Označme

$$B_n(\mathbf{X}) = \{\theta \in \Theta : S_n(\mathbf{X}, \theta) \notin C_\theta(\alpha)\}.$$

Potom pro všechna $F_X \in \mathcal{F}$ a tomuto rozdělení odpovídají parametr $\theta_X = t(F_X)$ platí

$$P_{F_X}[B_n(\mathbf{X}) \ni \theta_X] = P_{F_X}[S_n(\mathbf{X}, \theta_X) \notin C_{\theta_X}(\alpha)] = 1 - \alpha,$$

což bylo dokázat. □

Tvrzení 4.2 říká, že umíme-li sestavit interval spolehlivosti pro parametr, můžeme jej ihned využít k testování hypotéz o tomto parametru. Naopak máme-li test, můžeme s jeho pomocí sestavit interval spolehlivosti. Tento krok je však pracnější, protože vyžaduje otestování všech možných hodnot parametru. Množina hodnot parametru, pro které nezamítáme hypotézu, pak dává požadované pokrytí pro skutečný parametr, ale nemusí nutně tvořit interval.

Příklad. Máme náhodný výběr X_1, \dots, X_n z normálního rozdělení $F_X = N(\theta_X, \sigma_X^2) \in \mathcal{F}^B = \{N(\theta, \sigma^2), \theta \in \mathbb{R}, \sigma^2 > 0\}$.

Předpokládejme, že máme spočtený interval spolehlivosti (3.4) pro střední hodnotu normálního rozdělení s neznámým rozptylem. Potom zamítáme nulovou hypotézu $H_0 : \theta_X = \theta_0$ proti alternativě $H_1 : \theta_X \neq \theta_0$, pokud

$$\theta_0 \notin \left(\bar{X}_n - t_{n-1} \left(1 - \frac{\alpha}{2}\right) \frac{S_n}{\sqrt{n}}, \bar{X}_n + t_{n-1} \left(1 - \frac{\alpha}{2}\right) \frac{S_n}{\sqrt{n}} \right).$$

Tj. interval spolehlivosti obsahuje ty hodnoty parametru, pro které bychom nezamítli nulovou hypotézu.

Na druhou stranu pokud pro test $H_0 : \theta_X = \theta$ proti alternativě $H_1 : \theta_X \neq \theta$ použijeme testovou statistiku

$$T_n(\theta) = \frac{\sqrt{n}(\bar{X}_n - \theta)}{S_n},$$

(viz příklad (B) na str. 81). Pak výše uvedený interval spolehlivosti můžeme odvodit jako

$$\begin{aligned} & \{\theta \in \mathbb{R} : \text{nezamítáme } H_0 : \theta_X = \theta \text{ proti } H_1 : \theta_X \neq \theta\} \\ &= \{\theta \in \mathbb{R} : |T_n(\theta)| < t_{n-1}(1 - \alpha/2)\} \\ &= \{\theta \in \mathbb{R} : \left| \frac{\sqrt{n}(\bar{X}_n - \theta)}{S_n} \right| < t_{n-1}(1 - \alpha/2)\}. \end{aligned}$$

Cvičení. Jak by v předchozím příkladě vypadal interval spolehlivosti, který bychom odvodili z jednostranného test, tj. při testování $H_0 : \theta_X \leq \theta_0$ proti alternativě $H_1 : \theta_X > \theta_0$.

5. JEDNOVÝBĚROVÉ A PÁROVÉ PROBLÉMY PRO KVANTITATIVNÍ DATA

V této kapitole uvažujeme náhodný výběr X_1, \dots, X_n kvantitativních veličin s distribuční funkcí F_X patřící do modelu \mathcal{F} . Zajímá nás parametr $\theta_X = t(F_X)$. Chceme testovat hypotézy o tomto parametru, případně pro něj sestavit intervalový odhad.

5.1. JEDNOVÝBĚROVÝ KOLMOGOROVŮV-SMIRNOVŮV TEST

Jednovýběrový Kolmogorovův-Smirnovův test* testuje shodu distribuční funkce dat s určitou pevně danou distribuční funkcí. Je to neparametrický test, protože nepředpokládá žádný parametrický model.

Model: $\mathcal{F} = \{\text{všechna spojitá rozdělení}\}$

Testovaný parametr: Celá distribuční funkce F_X

Hypotéza a alternativa:

$$H_0 : F_X(x) = F_0(x) \quad \forall x \in \mathbb{R}, \quad H_1 : \exists x \in \mathbb{R} : F_X(x) \neq F_0(x),$$

kde F_0 je nějaká pevně specifikovaná spojitá distribuční funkce (bez neznámých parametrů).

Testová statistika je založena na empirické distribuční funkci \widehat{F}_n , s níž jsme se seznámili v kapitole 3.5.1 (viz str. 57). Její vlastnosti shrnuje věta 3.3. Empirická distribuční funkce je nestranným a konsistentním odhadem skutečné distribuční funkce v každém bodě. Navíc podle věty 3.3, bod (v), splňuje stejnoměrnou konsistenci, tj.

$$\sup_{x \in \mathbb{R}} |\widehat{F}_n(x) - F_X(x)| \xrightarrow[n \rightarrow \infty]{P} 0.$$

Testová statistika přebírá tuto supremální normu a zachycuje s ní největší celkový rozdíl mezi $\widehat{F}_n(x)$ a $F_0(x)$.

Testová statistika:

$$K_n = \sup_{x \in \mathbb{R}} |\widehat{F}_n(x) - F_0(x)|$$

Pokud hypotéza platí a F_0 je skutečná distribuční funkce dat, hodnota testové statistiky K_n bude blízko nuly. Hypotézu zamítneme, pokud se empirická distribuční funkce příliš liší od F_0 , tj. pokud je testová statistika příliš velká.

* Angl. *one-sample Kolmogorov-Smirnov test*

Označme

$$K_n^+ = \sup_{x \in \mathbb{R}} (\widehat{F}_n(x) - F_0(x)) \quad \text{a} \quad K_n^- = \sup_{x \in \mathbb{R}} (F_0(x) - \widehat{F}_n(x)).$$

Pak $K_n = \max(K_n^+, K_n^-)$.

Lemma 5.1 Je-li F_0 **spojitá**, tak platí

$$K_n^+ = \max_{1 \leq i \leq n} \left(\frac{i}{n} - F_0(X_{(i)}) \right), \quad K_n^- = \max_{1 \leq i \leq n} \left(F_0(X_{(i)}) - \frac{i-1}{n} \right).$$

Důkaz. Definujme si $X_{(0)} = -\infty$ a $X_{(n+1)} = +\infty$. Potom

$$\widehat{F}_n(x) = \frac{i}{n}, \quad \text{pro } x \in \langle X_{(i)}, X_{(i+1)} \rangle, \quad i = 0, 1, \dots, n.$$

Tedy s využitím výše uvedeného

$$\begin{aligned} K_n^+ &= \sup_{x \in \mathbb{R}} (\widehat{F}_n(x) - F_0(x)) = \max_{0 \leq i \leq n} \sup_{X_{(i)} \leq x < X_{(i+1)}} (\widehat{F}_n(x) - F_0(x)) \\ &= \max_{0 \leq i \leq n} \left(\frac{i}{n} - \inf_{X_{(i)} \leq x < X_{(i+1)}} F_0(x) \right) \\ &= \max_{0 \leq i \leq n} \left(\frac{i}{n} - F_0(X_{(i)}) \right) = \max_{1 \leq i \leq n} \left(\frac{i}{n} - F_0(X_{(i)}) \right), \end{aligned}$$

kde v poslední rovnosti jsme využili toho, že $F_0(X_{(0)}) = 0$ a že $1 - F_0(X_{(n)}) \geq 0$.

Podobně se dá upravit výraz pro K_n^- :

$$\begin{aligned} K_n^- &= \sup_{x \in \mathbb{R}} (F_0(x) - \widehat{F}_n(x)) = \max_{0 \leq i \leq n} \sup_{X_{(i)} \leq x < X_{(i+1)}} (F_0(x) - \widehat{F}_n(x)) \\ &= \max_{0 \leq i \leq n} (F_0(X_{(i+1)}) - \frac{i}{n}) = \max_{0 \leq i \leq n-1} (F_0(X_{(i+1)}) - \frac{i}{n}) \\ &= \max_{1 \leq i \leq n} (F_0(X_{(i)}) - \frac{i-1}{n}), \end{aligned}$$

kde jsme v předposlední rovnosti využili toho, že $F_0(X_{(n+1)}) = 1$ a že $F_0(X_{(1)}) \geq 0$.
V poslední rovnosti jsme pak pouze posunuli indexy. □

Poznámka. Předchozí lemma má několik důležitých důsledků.

- Testová statistika K_n se počítá pomocí Lemmatu 5.1, nikoli podle její definice. K jejímu výpočtu není třeba znát \widehat{F}_n .
- Platí-li hypotéza, $F_0(X_{(i)})$ má podle věty 2.13 beta rozdělení, jehož parametry nezávisí na F_0 . Dá se ukázat, že rozdělení K_n za platnosti hypotézy také nezávisí na F_0 (tj. K_n je pivotální).
- Z lemmatu 5.1 lze odvodit přesné rozdělení testové statistiky za platnosti hypotézy. Jedná se ovšem o netriviální výpočet, který je navíc i numericky obtížný. Proto se přesné rozdělení K_n zpravidla používá jen při malých rozsazích výběru n , pro které jsou kritické hodnoty tabelovány.

Asymptotické rozdělení testové statistiky za platnosti hypotézy je určeno následujícím tvrzením, které rozšiřuje výsledek uvedený ve větě 3.3, bod (v).

Tvrzení 5.2 Necht' X_1, \dots, X_n je náhodný výběr ze **spojitého** rozdělení s distribuční funkcí F_X . Potom

$$\sqrt{n} \sup_{x \in \mathbb{R}} |\widehat{F}_n(x) - F_X(x)| \xrightarrow[n \rightarrow \infty]{d} Z,$$

kde náhodná veličina Z má distribuční funkci danou předpisem*

$$G(y) = \begin{cases} 1 - 2 \sum_{k=1}^{\infty} (-1)^{k+1} e^{-2k^2 y^2}, & y > 0, \\ 0, & y \leq 0. \end{cases} \quad (5.1)$$

Distribuční funkce $G(y)$ určuje limitní rozdělení normalizované testové statistiky $\sqrt{n}K_n$ za platnosti hypotézy, tj. pro $F_X = F_0$. Toto rozdělení není normální, jak jsme byli doposud u limitních rozdělení zvyklí. Za povšimnutí stojí, že toto rozdělení nezávisí na volbě F_0 (tj. je pivotální). Důkaz tvrzení 5.2 náleží do pokročilé teorie pravděpodobnosti, my jej neuvádíme.

Nyní již můžeme určit kritickou hodnotu pro zamítání H_0 , aby měl test asymptotickou hladinu α . Označme α -kvantil rozdělení s distribuční funkcí G symbolem $k_\alpha = G^{-1}(\alpha)$. Hypotézu budeme zamítat, pokud $\sqrt{n}K_n$ překročí $k_{1-\alpha}$.

Kritický obor:

$$H_0 \text{ zamítneme} \Leftrightarrow \sqrt{n}K_n \geq k_{1-\alpha}. \quad (5.2)$$

Díky tvrzení 5.2 víme, že tento test má asymptoticky hladinu α .

P-hodnota: $p = 1 - G(\sqrt{n}k_n)$, kde k_n je napozorovaná hodnota statistiky K_n . Jedná se zde o asymptotickou p-hodnotu.

Poznámka.

- Všimněme si, že za alternativy

$$K_n \xrightarrow[n \rightarrow \infty]{P} \sup_{x \in \mathbb{R}} |F_X(x) - F_0(x)| > 0$$

a odtud se dá odvodit konzistence testu. Výhodou Kolmogorovova-Smirnovova testu je tedy jeho universalita (reaguje na jakýkoli rozdíl v rozdělení dat proti hypotéze) a absence předpokladů o rozdělení F_X .

- Tento test má relativně malou sílu proti konkrétnímu typu porušení H_0 (např. změna střední hodnoty). Pokud tušíme, jaké porušení H_0 je pro danou aplikaci neočekávanější nebo nejvíce relevantní, je lepší použít test, který je zaměřen na tento typ porušení H_0 .
- Tento test lze zformulovat i jako jednostranný proti alternativě $H_1' : F_X(x) \geq F_0(x), \exists x \in \mathbb{R} : F_X(x) > F_0(x)$ nebo $H_1'' : F_X(x) \leq F_0(x), \exists x \in \mathbb{R} : F_X(x) < F_0(x)$. Jako testovou statistiku pak použijeme buď K_n^+ anebo K_n^- a zamítáme pro jejich velké hodnoty. Pro určení (asymptotických) kritických hodnot však nelze využít tvrzení 5.2 a muselo by se hledat asymptotické rozdělení statistiky $\sqrt{n}K_n^+$ (resp. $\sqrt{n}K_n^-$).

* Viz např. kapitola 2.1.5 knihy [Serfling \(1980\)](#).

INTERVALY SPOLEHLIVOSTI PRO F_X

Obraťme nyní pozornost k problému sestrojení intervalu spolehlivosti pro distribuční funkci. Jestliže máme dané pevné $x \in S_X = \{y : F_X(y) \in (0, 1)\}$ a chceme intervalový odhad pouze pro hodnotu $F_X(x)$, můžeme vyjít z věty 3.3, bod (iii), a použít postup uvedený v příkladě na str. 56 v kapitole 3.4.2. Dostaneme interval

$$IS_n(x) = \left(\widehat{F}_n(x) - \frac{u_{1-\frac{\alpha}{2}} \sqrt{\widehat{F}_n(x)(1-\widehat{F}_n(x))}}{\sqrt{n}}, \widehat{F}_n(x) + \frac{u_{1-\frac{\alpha}{2}} \sqrt{\widehat{F}_n(x)(1-\widehat{F}_n(x))}}{\sqrt{n}} \right).$$

Pro tento interval platí

$$P[IS_n(x) \ni F_X(x)] \xrightarrow{n \rightarrow \infty} 1 - \alpha, \forall x \in S_X$$

a mluvíme o něm jako o „bodovém“ intervalu spolehlivosti* pro $F_X(x)$.

Co když ale nemáme předem dané x , nýbrž chceme interval, který by pokryl hodnotu distribuční funkce kdekoli, třeba i v mnoha bodech zároveň? K tomu nemůžeme použít postup uvedený výše, ale znovu využijeme tvrzení 5.2. Máme totiž

$$P\left[\sqrt{n}|\widehat{F}_n(x) - F_X(x)| < k_{1-\alpha}, \forall x \in \mathbb{R}\right] = P\left[\sqrt{n} \sup_{x \in \mathbb{R}} |\widehat{F}_n(x) - F_X(x)| < k_{1-\alpha}\right] \xrightarrow{n \rightarrow \infty} 1 - \alpha.$$

Sestavíme-li tedy pro každé $x \in \mathbb{R}$ interval

$$B_n(x) = \left(\widehat{F}_n(x) - \frac{k_{1-\alpha}}{\sqrt{n}}, \widehat{F}_n(x) + \frac{k_{1-\alpha}}{\sqrt{n}} \right),$$

potom

$$P[B_n(x) \ni F_X(x), \forall x \in S_X] \xrightarrow{n \rightarrow \infty} 1 - \alpha.$$

Intervalům vytvářejícím oblast, v níž se se zadanou pravděpodobností nachází celý průběh nějaké neznámé funkce, se říká *pás spolehlivosti*[†]. Protože hranice pásu spolehlivosti pro distribuční funkci založené na Kolmogorovově-Smirnovově statistice mohou ležet mimo přirozený rozsah $\langle 0, 1 \rangle$, předefinujeme dolní mez na $\max\{0, \widehat{F}_n(x) - k_{1-\alpha}/\sqrt{n}\}$ a horní mez na $\min\{1, \widehat{F}_n(x) + k_{1-\alpha}/\sqrt{n}\}$.[‡]

PORUŠENÍ PŘEDPOKLADŮ TESTU

Rozdělení F_0 není spojité V tomto případě lze použít statistiku K_n , neplatí však pro ni tvrzení 5.2. Pokud bychom tento fakt ignorovali a použili pro vyhodnocení testu kvantil $k_{1-\alpha}$, tak výsledný test bude (asymptoticky) konzervativní a tím pádem bude mít i menší sílu. Je třeba také dát pozor, že v tomto případě nelze testovou statistiku K_n počítat pomocí lemmatu 5.1.

* Angl. *pointwise confidence interval* † Angl. *confidence bands* ‡ Existuje samozřejmě řada jiných způsobů, jak sestavit pás spolehlivosti pro distribuční funkci.

Rozdělení F_0 je sice spojité, ale v datech jsou shody. Striktně vzato, pokud data pochází ze spojitého rozdělení, tak je nulová pravděpodobnost, že bychom měli nějaká pozorování se stejnými hodnotami (neboli shody*). V aplikacích však typicky vznikají shody kvůli zaokrouhlování. Tedy formálně pozorujeme vlastně $\tilde{X}_1, \dots, \tilde{X}_n$, kde \tilde{X}_i je zaokrouhlená náhodná veličina X_i . Čistě z teoretického hlediska tedy za nulové hypotézy empirická distribuční funkce

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{\tilde{X}_i \leq x\}$$

odhaduje distribuční funkci $\tilde{F}_0(x) = P(\tilde{X}_i \leq x)$ zaokrouhlené náhodné veličiny. Nicméně test se dá nadále používat jako přibližný test, pokud \tilde{F}_0 není příliš odlišné od F_0 . Přesněji pokud se dá očekávat, že

$$\sqrt{n} \sup_{x \in \mathbb{R}} |\tilde{F}_0(x) - F_0(x)|,$$

není příliš velké, což je často v aplikacích splněno.

Hypotéza není jednoduchá. Všimněme si, že F_0 musí být známa přesně (tj. nesmí obsahovat neznámé parametry ani jejich odhady). Předpokládejme, že potřebujeme testovat hypotézy

$$H_0 : F_X \in \mathcal{F}_0, \quad H_1 : F_X \notin \mathcal{F}_0,$$

kde $\mathcal{F}_0 = \{F(x; \theta), \theta \in \Theta\}$ je nějaká parametrická rodina rozdělení (např. $\{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$). Potom je přirozené uvažovat jako testovou statistiku

$$\tilde{K}_n = \sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x; \hat{\theta}_n)|,$$

kde $\hat{\theta}_n$ je odhad skutečné hodnoty parametru θ_X . Je však nutné si uvědomit, že pro statistiku \tilde{K}_n již **neplatí** tvrzení 5.2. Navíc se ukazuje, že i za nulové hypotézy je asymptotické rozdělení \tilde{K}_n velmi komplikované a závisí na neznámém parametru θ_X . Pokud bychom ignorovali tento fakt a použili pro vyhodnocení testu kvantil $k_{1-\alpha}$, tak výsledný test bude silně konzervativní a tím pádem bude mít malou sílu.

Všechna výše zmíněna porušení předpokladů se dají řešit pomocí tzv. parametrického bootstrapu (viz přednáška *NMST545 Matematická statistika 4*).

Cvičení. Uvažujte, test s kritickým oborem

$$H_0 \text{ zamítneme} \Leftrightarrow \sqrt{n}K_n \leq k_{\alpha/2} \text{ nebo } \sqrt{n}K_n \geq k_{1-\alpha/2}.$$

1. Dodržujte tento test hladinu α (přesně nebo asymptoticky)?
2. Jak by se spočetla p-hodnota tohoto testu?
3. Je tento test konsistentní?
4. V čem by tento test mohl být lepší/horší než test s kritickým oborem (5.2)?

* Angl. *ties*

5.2. JEDNOVÝBĚROVÝ t -TEST

Jednovýběrový t -test* porovnává **střední hodnotu** dat s nějakou zvolenou konstantou. Tímto testem jsme se podrobně zabývali v Příkladě B na str. 81 a také v Příkladě C na str. 83. Rozdíl byl pouze v tom, jaký model jsme předpokládali.

Model: $\mathcal{F}^B = \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$ nebo $\mathcal{F}^C = \mathcal{L}_+^2$

Testovaný parametr: Střední hodnota $\mu_X = E X_i$

Hypotéza a alternativa:

$$H_0 : \mu_X = \mu_0, \quad H_1 : \mu_X \neq \mu_0,$$

kde μ_0 je předem daná konstanta.

Testová statistika:

$$T_n = \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{S_n},$$

kde \bar{X}_n je aritmetický průměr a S_n^2 je výběrový rozptyl.

Rozdělení testové statistiky za H_0 :

V modelu \mathcal{F}^B : $T_n \sim t_{n-1}$ (viz věta 2.10), v modelu \mathcal{F}^C : $T_n \stackrel{\text{as.}}{\sim} N(0, 1)$ (viz věta 2.9).

Tedy test bude **přesný**, pokud platí „menší“ model \mathcal{F}^B . Pro „větší“ model \mathcal{F}^C je tento test **asymptotický**. Podobně to platí dále pro p -hodnotu a intervalový odhad. V modelu \mathcal{F}^B budou p -hodnota a intervalový odhad přesné, v modelu \mathcal{F}^C pouze asymptotické.

Kritický obor:

$$H_0 \text{ zamítneme} \Leftrightarrow |T_n| \geq t_{n-1}(1 - \alpha/2),$$

kde $t_{n-1}(1 - \alpha/2)$ je $(1 - \alpha/2)$ -tý kvantil t -rozdělení s $n - 1$ stupni volnosti.

p -hodnota: $p = 2(1 - F_n(|t|))$, kde t je pozorovaná hodnota testové statistiky T_n a F_n je distribuční funkce rozdělení t_{n-1} .

Interval spolehlivosti pro μ_X : Interval spolehlivosti pro střední hodnotu normálního rozdělení je dán krajními body

$$\left(\bar{X}_n - t_{n-1}(1 - \frac{\alpha}{2}) \frac{S_n}{\sqrt{n}}, \bar{X}_n + t_{n-1}(1 - \frac{\alpha}{2}) \frac{S_n}{\sqrt{n}} \right).$$

Viz vzorec na str. 55 a následující příklad. Tento interval je v modelu \mathcal{F}^B přesný v modelu \mathcal{F}^C asymptotický.

Poznámka. T -test nutně nepotřebuje předpoklad normálního rozdělení, funguje jako asymptotický test pro libovolné rozdělení s konečným rozptylem. Pouze je potřeba mít k dispozici dostatek pozorování.

* Angl. *one-sample t-test*

Poznámka. T-test lze převést na jednostranný test:

$$\text{zamítáme } H'_0 : \mu_X \leq \mu_0 \text{ proti } H'_1 : \mu_X > \mu_0 \Leftrightarrow T_n \geq t_{n-1}(1 - \alpha).$$

Podobně zamítáme $H''_0 : \mu_X \geq \mu_0$ proti $H''_1 : \mu_X < \mu_0$, pokud testová statistika T_n je nižší než kritická hodnota $-t_{n-1}(1 - \alpha)$.

Všimněte si, že **nelze** tvrdit, že za nulové hypotézy H'_0 v modelu \mathcal{F}^B platí $T_n \sim t_{n-1}$, resp. $T_n \stackrel{\text{as.}}{\sim} N(0, 1)$ v modelu \mathcal{F}^C . Rozmyslete si, proč to však nebrání výše popsanému provedení testu.

5.3. JEDNOVÝBĚROVÝ ZNAMÉNKOVÝ TEST

Jednovýběrový znaménkový test* porovnává **medián** dat s pevně danou hodnotou. Je to neparametrický test, funguje pro jakékoli spojitě rozdělení.

Model: $\mathcal{F} = \{\text{všechna spojitá rozdělení}\}$

Testovaný parametr: Medián $m_X = F_X^{-1}(0.5)$

Hypotéza a alternativa:

$$H_0 : m_X = m_0, \quad H_1 : m_X \neq m_0,$$

kde m_0 je předem daná konstanta.

Testová statistika:

$$B_n = \sum_{i=1}^n \mathbb{1}\{X_i > m_0\}$$

(počet pozorování větších než m_0).

Věta 5.3 Nechť X_1, \dots, X_n je náhodný výběr z libovolného **spojitého** rozdělení s mediánem m_X . Pak

(i)

$$\sum_{i=1}^n \mathbb{1}\{X_i > m_X\} \sim \text{Bi}(n, \frac{1}{2}),$$

(ii)

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \left[\mathbb{1}\{X_i > m_X\} - \frac{1}{2} \right] \xrightarrow[n \rightarrow \infty]{d} N(0, \frac{1}{4}).$$

Poznámka. Věta 5.3 plyne z věty 2.3, části (iii) a (iv).

Přesné rozdělení testové statistiky za H_0 :

$$B_n \sim \text{Bi}(n, \frac{1}{2}), \quad (\text{viz věta 5.3(i)})$$

* Angl. *one-sample sign test*

Kritický obor (přesný test): Hypotézu budeme zamítat pro příliš malé nebo příliš velké hodnoty B_n .

$$H_0 \text{ zamítneme} \Leftrightarrow B_n \leq c_L(\alpha) \text{ nebo } B_n \geq c_U(\alpha),$$

kde

$$c_L(\alpha) = \max \left\{ k_1 \in \mathbb{N}_0 : P\left(\text{Bi}\left(n, \frac{1}{2}\right) \leq k_1\right) \leq \frac{\alpha}{2} \right\} = \max \left\{ k_1 \in \mathbb{N}_0 : \frac{1}{2^n} \sum_{j=0}^{k_1} \binom{n}{j} \leq \frac{\alpha}{2} \right\},$$

$$c_U(\alpha) = \min \left\{ k_2 \in \mathbb{N}_0 : P\left(\text{Bi}\left(n, \frac{1}{2}\right) \geq k_2\right) \leq \frac{\alpha}{2} \right\} = \min \left\{ k_2 \in \mathbb{N}_0 : \frac{1}{2^n} \sum_{j=k_2}^n \binom{n}{j} \leq \frac{\alpha}{2} \right\}.$$

Ze symetrie binomického rozdělení pro $p = \frac{1}{2}$ plyne, že $c_L(\alpha) + c_U(\alpha) = n$. Tento test má hladinu nejvýše α (přesné hladiny α nemusí být možné dosáhnout).

P-hodnota (přesná):

$$p = 2 \min \left\{ P\left(\text{Bi}\left(n, \frac{1}{2}\right) \leq y_n\right), P\left(\text{Bi}\left(n, \frac{1}{2}\right) \geq y_n\right) \right\} = 2 \min \{1 - G_0(y_n - 1), G_0(y_n)\},$$

kde G_0 je distribuční funkce $\text{Bi}(n, \frac{1}{2})$ a y_n je napozorovaná hodnota B_n .

Asymptotické rozdělení testové statistiky za H_0 :

$$Z_n = \frac{B_n - \frac{n}{2}}{\sqrt{\frac{n}{4}}} = \frac{2}{\sqrt{n}} \left(B_n - \frac{n}{2} \right) \stackrel{\text{as.}}{\sim} N(0, 1), \quad (\text{viz věta 5.3(ii)}).$$

Kritický obor (asymptotický test): Hypotézu budeme zamítat pro příliš malé nebo příliš velké hodnoty B_n .

$$H_0 \text{ zamítneme} \Leftrightarrow |Z_n| \geq u_{1-\alpha/2}.$$

P-hodnota (asymptotická): $p = 2(1 - \Phi(|z_n|))$, kde z_n je napozorovaná hodnota testové statistiky Z_n .

Interval spolehlivosti pro m_X : Viz intervaly spolehlivosti pro kvantily (kapitola 3.5.4).

Poznámka.

- K výpočtu testové statistiky vlastně nepotřebujeme znát konkrétní hodnoty X_i . Stačí nám jen vědět, kolik z nich překročilo hodnotu m_0 .
- Tento test lze převést na jednostranný test $H'_0 : m_X \geq m_0$ (nebo $\leq m_0$).
- Test lze snadno modifikovat na test o libovolném kvantilu, tj. na test hypotéz

$$H_0 : u_X(\beta) = u_0, \quad H_1 : u_X(\beta) \neq u_0,$$

kde $\beta \in (0, 1)$. Testová statistika $B_n = \sum_{i=1}^n \mathbb{1}\{X_i > u_0\}$ potom bude mít za nulové hypotézy rozdělení $\text{Bi}(n, 1 - \beta)$. Testování o kvantilu tedy převedeme na testování hodnoty parametru binomického rozdělení, čímž se budeme podrobně zabývat v kapitole 7.1.

Cvičení. Ukažte, že znaménkový test je consistentní.

Návod: Je jednodušší pracovat s asymptotickou verzí znaménkového testu.

PORUŠENÍ PŘEDPOKLADŮ

I když se zpravidla vyžaduje spojitost rozdělení F_X , tak pro dodržení (přesné či asymptotické) hladiny stačí, že za nulové hypotézy $P[X_i = m_0] = 0$. V aplikacích se však může stát, že i když se dá předpokládat, že sledovaná veličina je spojitá, tak vlivem zaokrouhlování jsou některá pozorování přesně rovna m_0 . Taková pozorování vyloučíme, protože nelze určit, zda před zaokrouhlením byla větší nebo menší než m_0 . Test pak provádíme na zmenšeném výběru.

5.4. JEDNOVÝBĚROVÝ WILCOXONŮV TEST

Jednovýběrový Wilcoxonův test* předpokládá **symetrické** rozdělení a porovnává **střed symetrie** s pevně danou konstantou.

Model: $\mathcal{F} = \{\text{spojitá rozdělení s hustotou } f \text{ splňující } \exists \delta \in \mathbb{R} : f(\delta - x) = f(\delta + x) \forall x \in \mathbb{R}\}$

Testovaný parametr: Střed symetrie δ_X

Poznámka. Model vyžaduje, aby **hustota** X_i byla **symetrická** kolem nějakého bodu δ_X . Pak musí platit $m_X = \delta_X$ a pokud $X_i \in \mathcal{L}^1$, pak i $E X_i = \mu_X = \delta_X$.

Hypotéza a alternativa:

$$H_0 : \delta_X = \delta_0, \quad H_1 : \delta_X \neq \delta_0,$$

kde δ_0 je předem daná konstanta.

Poznámka. Za platnosti modelu \mathcal{F} je hypotéza H_0 ekvivalentní hypotéze $H_0^* : m_X = \delta_0$ (test na medián). Pokud navíc $X_i \in \mathcal{L}^1$, pak je hypotéza H_0 též ekvivalentní hypotéze $H_0^{**} : \mu_X = \delta_0$ (test na střední hodnotu).

Testová statistika: Nechť $Z_i \stackrel{\text{df}}{=} X_i - \delta_0$. Definujme

$$W_n = \sum_{i \in \mathcal{I}} R_i,$$

kde $\mathcal{I} = \{i \in \{1, \dots, n\} : Z_i > 0\}$ je množina všech indexů takových, že Z_i má kladné znaménko a R_1, R_2, \dots, R_n jsou pořadí absolutních hodnot $|Z_i|$ mezi všemi absolutními hodnotami $|Z_1|, \dots, |Z_n|$.

Poznámka. Testová statistika W_n jednovýběrového Wilcoxonova testu může nabývat hodnot z množiny $\{0, 1, \dots, \frac{n(n+1)}{2}\}$. Spočítá se následujícím způsobem:

1. Spočítáme odchylky $Z_i = X_i - \delta_0$ a určíme množinu indexů \mathcal{I} .
2. Spočteme $|Z_1|, \dots, |Z_n|$.

* Angl. *one-sample Wilcoxon test, Wilcoxon signed rank test*

3. Seřadíme všechny $|Z_i|$ od nejmenší do největší a získáme uspořádaný výběr

$$0 < |Z|_{(1)} < |Z|_{(2)} < \dots < |Z|_{(n)}.$$

4. Určíme pořadí R_i náhodné veličiny $|Z_i|$ mezi všemi $|Z|_{(1)}, \dots, |Z|_{(n)}$. Platí $|Z_i| = |Z|_{(R_i)}$.

5. Sečteme pořadí R_i pro $i \in \mathcal{I}$.

Velikost množiny \mathcal{I} je rovna počtu pozorování, pro něž platí $X_i > \delta_0$ (srv. s testovou statistikou znaménkového testu).

Tvrzení 5.4 Necht' X_1, \dots, X_n je náhodný výběr z libovolného **spojitého** rozdělení splňujícího model \mathcal{F} a necht' **platí hypotéza** $H_0 : \delta_X = \delta_0$. Pak

(i)

$$E_{H_0} W_n = \frac{n(n+1)}{4}, \quad \text{var}_{H_0}(W_n) = \frac{n(n+1)(2n+1)}{24};$$

(ii)

$$\frac{W_n - E_{H_0} W_n}{\sqrt{\text{var}_{H_0}(W_n)}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Důkaz. Bez újmy na obecnosti uvažujme $\delta_0 = 0$. Zaveďme náhodné veličiny $\Delta_i = \text{sign}(Z_i)$. Všimněme si, že

$$W_n = \sum_{i=1}^n R_i \mathbb{1}\{\Delta_i = 1\}.$$

Náhodné veličiny $\Delta_1, \dots, \Delta_n$ jsou za H_0 nezávislé stejně rozdělené takové, že

$$P(\Delta_i = 1) = P(\Delta_i = -1) = \frac{1}{2}.$$

Odtud snadno spočteme

$$E \Delta_i = 0, \quad E \Delta_i^2 = 1.$$

Důkaz si rozdělíme do tří kroků.

1. Ukážeme, že $(R_1, \dots, R_n)^T$ a $(\Delta_1, \dots, \Delta_n)^T$ jsou nezávislé.

Nejdříve si všimněme, že náhodný vektor $(R_1, \dots, R_n)^T$ je funkcí náhodný vektoru $(|Z_1|, \dots, |Z_n|)^T$. Tudíž stačí dokázat, že náhodné vektory $(|Z_1|, \dots, |Z_n|)^T$ a $(\Delta_1, \dots, \Delta_n)^T$ jsou nezávislé.

Za tímto účelem si dále všimněme, že náhodné vektory $(\binom{|Z_1|}{\Delta_1}), \dots, (\binom{|Z_n|}{\Delta_n})$ jsou nezávislé. Tedy zbývá pouze dokázat nezávislost $|Z_i|$ a Δ_i .

Pro $\forall z > 0$ platí

$$\begin{aligned} P[|Z_i| \leq z, \Delta_i = 1] &= P[0 \leq Z_i \leq z] = \frac{1}{2} P[-z \leq Z_i \leq z] \\ &= \frac{1}{2} P[0 \leq |Z_i| \leq z] = P[\Delta_i = 1] P[|Z_i| \leq z], \end{aligned}$$

kde jsme v druhé rovnosti využili toho, že rozdělení Z_i je (za nulové hypotézy) symetrické kolem nuly. Tedy $|Z_i|$ a Δ_i jsou vskutku nezávislé.

2. Vyjádříme si W_n pomocí R_i a Δ_i .

Máme

$$\begin{aligned} \sum_{i=1}^n R_i \mathbb{1}\{\Delta_i = 1\} + \sum_{i=1}^n R_i \mathbb{1}\{\Delta_i = -1\} &= \sum_{i=1}^n R_i = \frac{n(n+1)}{2}, \\ \sum_{i=1}^n R_i \mathbb{1}\{\Delta_i = 1\} - \sum_{i=1}^n R_i \mathbb{1}\{\Delta_i = -1\} &= \sum_{i=1}^n R_i \Delta_i. \end{aligned}$$

„Zprůměrováním“ výše uvedeným rovnic a využitím toho, že $W_n = \sum_{i=1}^n R_i \mathbb{1}\{\Delta_i = 1\}$ tedy dostaneme

$$W_n = \frac{n(n+1)}{4} + \frac{1}{2} \sum_{i=1}^n R_i \Delta_i. \quad (5.3)$$

3. Výpočet $E_{H_0} W_n$ a $\text{var}_{H_0}(W_n)$.

Ze vzorce (5.3) nyní s využitím nezávislosti R_i a Δ_i a toho, že $E \Delta_i = 0$ plyne

$$E_{H_0} W_n = \frac{n(n+1)}{4} + \frac{1}{2} \sum_{i=1}^n E R_i E \Delta_i = \frac{n(n+1)}{4}.$$

Dále

$$\text{var}_{H_0}(W_n) = \frac{1}{4} \text{var} \left(\sum_{i=1}^n R_i \Delta_i \right) = \frac{1}{4} \sum_{i=1}^n \text{var} (R_i \Delta_i) + \frac{1}{4} \sum_{i=1}^n \sum_{j=1, j \neq i}^n \text{cov} (R_i \Delta_i, R_j \Delta_j).$$

Spočítejme si tedy

$$\text{var} (R_i \Delta_i) = E (R_i \Delta_i)^2 = E R_i^2 E \Delta_i^2 = \frac{1}{n} \sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{6n} = \frac{(n+1)(2n+1)}{6},$$

kde jsme využili $E R_i \Delta_i = 0$, $E \Delta_i^2 = 1$ a věty 2.16(i), dle které $P[R_i = k] = \frac{1}{n}$ pro všechna $i, k \in \{1, \dots, n\}$.

Dále pro $i \neq j$ počítejme

$$\text{cov} (R_i \Delta_i, R_j \Delta_j) = E (R_i \Delta_i R_j \Delta_j) = E (R_i R_j) E \Delta_i E \Delta_j = 0,$$

kde jsme využili nezávislosti R_i a Δ_i .

Tedy celkem dostáváme

$$\text{var}_{H_0}(W_n) = \frac{1}{4} \sum_{i=1}^n \frac{(n+1)(2n+1)}{6} = \frac{n(n+1)(2n+1)}{24}.$$

□

Poznámka.

- Důkaz asymptotické normality vynecháváme. Důkaz je obtížný v tom, že pořadí R_1, \dots, R_n nejsou nezávislé náhodné veličiny.
- Hypotézu budeme zamítat pro příliš malé nebo příliš velké hodnoty W_n .
- Není-li n příliš velké, lze **za nulové hypotézy** nalézt i přesné rozdělení testové statistiky W_n . To se dá nahlédnout následovně. Z důkazu tvrzení 5.4 víme, že testová statistika se dá napsat jako funkce sdruženého rozdělení vektoru pořadí $\mathbf{R} = (R_1, \dots, R_n)^\top$ a vektoru znamének $(\Delta_1, \dots, \Delta_n)^\top$, přičemž tyto náhodné vektory jsou za H_0 nezávislé. Dále rozdělení \mathbf{R} za nulové hypotézy nám dává věta 2.15. Náhodné veličiny $\Delta_1, \dots, \Delta_n$ jsou za H_0 nezávislé stejně rozdělené takové, že $P(\Delta_i = 1) = P(\Delta_i = -1) = \frac{1}{2}$. Tudíž máme všechny potřebné informace k tomu, abychom odvodili přesné rozdělení W_n za nulové hypotézy.

Zatímco pro malé hodnoty rozsahů výběru n jsou kritické hodnoty tabelovány, pro střední a velké hodnoty n se využívá asymptotická normalita testové statistiky W_n za nulové hypotézy (Tvrzení 5.4(ii)).

Asymptotické rozdělení testové statistiky za H_0 :

$$U_n = \frac{W_n - E_{H_0} W_n}{\sqrt{\text{var}_{H_0}(W_n)}} = \frac{W_n - \frac{n(n+1)}{4}}{\sqrt{\frac{n(n+1)(2n+1)}{24}}} \stackrel{\text{as.}}{\approx} N(0, 1)$$

Kritický obor (asymptotický test):

$$H_0 \text{ zamítneme} \Leftrightarrow |U_n| \geq u_{1-\alpha/2}.$$

P-hodnota (asymptotická): $p = 2(1 - \Phi(|u_n|))$, kde u_n je napozorovaná hodnota testové statistiky U_n .

Poznámka. Jednovýběrový Wilcoxonův test bere v úvahu i velikost odchylek od δ_0 , nikoli jen jejich znaménko (jako znaménkový test). Jeho síla pro testování mediánu je obecně větší než síla znaménkového testu. Hladinu však dodržuje pouze tehdy, je-li rozdělení jednotlivých pozorování symetrické, zatímco znaménkový test takový předpoklad nevyžaduje.

PORUŠENÍ PŘEDPOKLADŮ

Shody kvůli zaokrouhlování. Kvůli zaokrouhlování bývají v aplikacích v datech často shody. V tomto případě jako u znaménkového testu nejdříve odstraníme pozorování, která se rovnají přesně δ_0 . Testová statistika W_n se pak spočte ze zbývajících dat a v případě shod pracuje s tzv. *průměrným pořadím*, viz. (2.6). Pak se dá ukázat, že za nulové hypotézy

$$\frac{W_n - \frac{n(n+1)}{4}}{\sqrt{\frac{n(n+1)(2n+1)}{24} - kor.}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1),$$

kde n je již (případně zmenšený) rozsah výběru a $kor.$ je korekce rozptylu daná předpisem*

$$kor. = \frac{1}{48} \sum_z (t_z^3 - t_z),$$

kde t_z značí počet, kolikrát se mezi hodnotami $|Z_1|, \dots, |Z_n|$ vyskytla hodnota z . Suma \sum_z pak značí sčítání přes všechny rozdílné hodnoty množiny $\{|Z_1|, \dots, |Z_n|\}$.

Za povšimnutí stojí, že bez úpravy jmenovatele pomocí $kor.$ by byl test konzervativní.

Nesymetrie. V případě, že hustota f není symetrická, pak testovaný parametr není medián pozorování X_i , ale tzv. *pseudo-medián*, což je medián náhodné veličiny $\frac{X_1+X_2}{2}$. Problém pseudo-mediánu je jeho obtížná interpretace. Obecně lze pouze říct, že jeho hodnota je někde mezi mediánem m_X a střední hodnotou $E X_i$ (pokud tato střední hodnota existuje).

Dalším nepříjemným důsledkem asymetrie pak je, že i pokud se díváme na jednovýběrový Wilcoxonův test jako na test o pseudo-mediánu, tak jeho skutečná hladina (ať již přesná nebo asymptotická) je odlišná od předepsaného α . Nicméně se ukazuje, že tato odchylka je vcelku malá i pro natolik asymetrická rozdělení jako je například exponenciální. Pokud tedy data nevykazují naprosto očividnou asymetrii, je tedy hlavním problémem interpretace pseudo-mediánu.

5.5. JEDNOVÝBĚROVÝ χ^2 -TEST NA ROZPTYL

Jednovýběrový χ^2 -test na rozptyl[†] je přesný test vyžadující normální rozdělení pozorovaných dat.

Model: $\mathcal{F} = \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$

Testovaný parametr: Rozptyl $\sigma_X^2 = \text{var } X_i$.

Hypotéza a alternativa:

$$H_0 : \sigma_X^2 = \sigma_0^2, \quad H_1 : \sigma_X^2 \neq \sigma_0^2,$$

kde σ_0^2 je předem daná konstanta.

Testová statistika:

$$\frac{(n-1)S_n^2}{\sigma_0^2},$$

kde S_n^2 je výběrový rozptyl (viz definice 2.4).

Přesné rozdělení testové statistiky za H_0 :

$$\frac{(n-1)S_n^2}{\sigma_0^2} \sim \chi_{n-1}^2 \quad (\text{viz věta 2.8}(i)).$$

* Viz např. Hollander et al. (2013), str. 42. † Angl. *one-sample chi-square variance test*

Kritický obor: Hypotézu zamítneme, pokud se výběrový rozptyl příliš liší od hypotetického rozptylu, tj. pokud je testová statistika buď moc velká nebo moc malá.

$$H_0 \text{ zamítneme} \Leftrightarrow \frac{(n-1)S_n^2}{\sigma_0^2} \leq \chi_{n-1}^2(\alpha/2) \text{ nebo } \frac{(n-1)S_n^2}{\sigma_0^2} \geq \chi_{n-1}^2(1-\alpha/2),$$

kde $\chi_{n-1}^2(\alpha/2)$ a $\chi_{n-1}^2(1-\alpha/2)$ jsou po řadě $(\alpha/2)$ -tý a $(1-\alpha/2)$ -tý kvantil χ^2 rozdělení s $n-1$ stupni volnosti.

P-hodnota: $p = 2 \min\{1 - G_{n-1}(s), G_{n-1}(s)\}$, kde s je pozorovaná hodnota testové statistiky a G_{n-1} je distribuční funkce rozdělení χ_{n-1}^2 .

Interval spolehlivosti pro σ_X^2 :

$$\left(\frac{(n-1)S_n^2}{\chi_{n-1}^2(1-\alpha/2)}, \frac{(n-1)S_n^2}{\chi_{n-1}^2(\alpha/2)} \right) \quad (\text{viz (3.3)}).$$

Cvičení. Ukažte, že jednovýběrový χ^2 -test na rozptyl je konsistentní.

Návod: Uvažujte pro jednoduchost jednostrannou hypotézu (a alternativu) a uvědomte si, že $\frac{\chi_{n-1}^2(\beta)}{n} \xrightarrow[n \rightarrow \infty]{} 1$ pro všechna $\beta \in (0, 1)$.

Poznámka.

- Při **porušení** předpokladu **normality** tento test **nedodrhuje hladinu** ani asymptoticky. V tomto případě lze zkonstruovat test na základě asymptotického rozdělení S_n^2 , viz věta 2.6(iii).
- Tento test lze převést na jednostranný test:

$$\text{zamítáme } H_0 : \sigma_X^2 \leq \sigma_0^2 \text{ proti } H_1' : \sigma_X^2 > \sigma_0^2 \Leftrightarrow \frac{(n-1)S_n^2}{\sigma_0^2} \geq \chi_{n-1}^2(1-\alpha).$$

Podobně hypotéza $H_0'' : \sigma_X^2 \geq \sigma_0^2$ proti alternativě $H_1'' : \sigma_X^2 < \sigma_0^2$, se zamítá pouze pro příliš malé hodnoty testové statistiky, přičemž kritická hodnota je $\chi_{n-1}^2(\alpha)$.

5.6. PÁROVÉ TESTY

Uvažujme náhodný výběr

$$\begin{pmatrix} X_1 \\ Y_1 \end{pmatrix}, \dots, \begin{pmatrix} X_n \\ Y_n \end{pmatrix}$$

dvousložkových náhodných vektorů s dvourozměrnou distribuční funkcí. Chceme porovnat nějakou charakteristiku marginálního rozdělení F_X náhodné veličiny X_i se stejnou charakteristikou marginálního rozdělení F_Y náhodné veličiny Y_i . Pozorování X_i a Y_i ovšem nejsou nezávislá.

Hlavní myšlenka párových testů je jednoduchá: Vezmeme rozdíly $Z_i = X_i - Y_i$ (jež tvoří náhodný výběr z nějakého jednorozměrného rozdělení) a na tyto rozdíly provedeme vhodný jednovýběrový test. Musíme se však zamyslet na tím, jestli hypotéza

testovaná jednovýběrovým testem provedeným na Z_i má nějakou rozumnou interpretaci pro porovnání rozdělení X_i a Y_i . Někdy tomu tak je, ale v řadě případů taková interpretace neexistuje (např. párový Kolmogorovův-Smirnovův test rozumnou interpretaci nemá).

Nechť například jednovýběrový test provedený na rozdíly Z_i testuje střední hodnotu, třeba $H_0 : E Z_i = 0$. Tato hypotéza je splněna právě tehdy, když $E X_i = E Y_i$ a výsledný test tedy testuje rovnost středních hodnot X_i a Y_i .

U jiných charakteristik toto neplatí: testujeme-li nulovost mediánu Z_i , neznamená to bez dalších předpokladů, že se za platnosti této hypotézy rovnají mediány X_i a Y_i . Testování rozptylu Z_i jednovýběrovým testem pak neříká vůbec nic o tom, jak a v čem se liší rozdělení X_i od rozdělení Y_i .

Párové testy lze použít pouze na intervalové a poměrové veličiny, jinak by rozdíly hodnot neměly smysluplnou interpretaci. Typicky je používáme na *uspořádané dvojice* měření téže veličiny na dvou přirozeně spárovaných jednotkách (např. levé oko – pravé oko, manžel – manželka) nebo dvě opakovaná měření téže veličiny na téže jednotce (např. před zásahem – po zásahu, loni – letos).

HYPOTÉZA NULOVÉHO EFEKTU

V aplikacích vyjadřuje náhodný vektor $(X_i, Y_i)^T$ měření před a po nějakém ošetření*. Nulová hypotéza pak říká, že ošetření nemělo žádný vliv. Tj. testujeme

$$H_0 : F_X(x) = F_Y(x), \forall x \in \mathbb{R} \quad H_1 : \exists x \in \mathbb{R} \quad F_X(x) \neq F_Y(x), \quad (5.4)$$

kde F_X a F_Y jsou (marginální) distribuční funkce náhodných veličin X_i a Y_i .

Je třeba si uvědomit, že každý z níže uvedených testů se zaměřuje pouze na určitý způsob porušení nulové hypotézy vyjádřené v (5.4).

5.7. PÁROVÝ t -TEST

Párový t -test[†] se provádí jako jednovýběrový t -test na rozdíly Z_i .

Model: $\mathcal{F}_n = \{Z_i = X_i - Y_i \sim N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$ nebo $\mathcal{F}_{as} = \{Z_i = X_i - Y_i \in \mathcal{L}_+^2\}$

Testované parametry: Střední hodnoty $\mu_X = E X_i$ a $\mu_Y = E Y_i$.

Hypotéza a alternativa:

$$H_0 : \mu_X - \mu_Y = \delta_0, \quad H_1 : \mu_X - \mu_Y \neq \delta_0,$$

kde δ_0 je předem daná konstanta (obvykle $\delta_0 = 0$).

Testová statistika:

$$T_n = \frac{\sqrt{n} (\bar{Z}_n - \delta_0)}{S_Z},$$

kde \bar{Z}_n je aritmetický průměr rozdílů Z_i (což je rovno $\bar{X}_n - \bar{Y}_n$) a S_Z je výběrová směrodatná odchylka rozdílů Z_i .

* Angl. *treatment* † Angl. *paired t-test*

Poznámka. Všimněme si, že

$$S_Z^2 = \frac{1}{n-1} \sum_{i=1}^n (Z_i - \bar{Z}_n)^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - Y_i - \bar{X}_n + \bar{Y}_n)^2 = S_X^2 - 2S_{X,Y} + S_Y^2,$$

kde S_X^2 a S_Y^2 jsou příslušné výběrové rozptyly a $S_{X,Y}$ je výběrová kovariance. Tj. testovou statistiku můžeme přepsat jako

$$T_n = \frac{\bar{X}_n - \bar{Y}_n - \delta_0}{\sqrt{S_X^2/n + S_Y^2/n - 2S_{X,Y}/n}},$$

což připomíná testovou statistiku dvouvýběrového t -testu v případě shodných rozsahů výběr (viz kapitola 6.2). V našem případě je však ve jmenovateli (pod odmocninou) navíc $-2S_{X,Y}/n$. Jelikož je zpravidla $S_{X,Y} > 0$ (neboť X_i a Y_i jsou typicky pozitivně korelované), tak mylné použití dvouvýběrového t -testu na párový problém by vedlo ke konzervativnímu testu s menší silou.

Rozdělení testové statistiky za H_0 :

$$\text{V modelu } \mathcal{F}_n : T_n \sim t_{n-1}, \quad \text{v modelu } \mathcal{F}_{as} : T_n \stackrel{as.}{\sim} N(0, 1).$$

Podobně jako u jednovýběrového t -testu (kapitola 5.3) bude tento test **přesný** v „menším“ modelu \mathcal{F}_n . Pro „větší“ model \mathcal{F}_{as} je pak tento test **asymptotický**. Podobně to platí dále pro p -hodnotu a intervalový odhad. V modelu \mathcal{F}_n budou p -hodnota a intervalový odhad přesné, v modelu \mathcal{F}_{as} pouze asymptotické.

Kritický obor:

$$H_0 \text{ zamítneme} \Leftrightarrow |T_n| \geq t_{n-1}(1 - \alpha/2),$$

kde $t_{n-1}(1 - \alpha/2)$ je $(1 - \alpha/2)$ -tý kvantil t -rozdělení s $n - 1$ stupni volnosti.

P -hodnota: $p = 2(1 - G_{n-1}(|t|))$, kde t je pozorovaná hodnota testové statistiky a G_{n-1} je distribuční funkce rozdělení t_{n-1} .

Interval spolehlivosti pro $\mu_X - \mu_Y$: Vypracujte samostatně.

Poznámka. Pro $\delta_0 = 0$ bychom se mohli na tento test dívat jako na test hypotézy nulového efektu (5.4). Z tohoto pohledu bude test citlivý **na rozdíl ve středních hodnotách** (tj. bude konzistentní, pokud budou střední hodnoty rozdílné). Naopak test nebude konzistentní, pokud sice H_0 v (5.4) platit nebude, ale bude $E Z_i = 0$. Tj. ošetření nebude mít efekt na střední hodnotu $E Y_i$, ale pouze například na rozptyl $\text{var } Y_i$.

5.8. PÁROVÝ ZNAMÉNKOVÝ TEST

Párový znaménkový test* se provádí jako jednovýběrový znaménkový test na rozdíly Z_i . Předpokládá se spojitost rozdílů Z_i .

* Angl. *paired sign test*

Model: $\mathcal{F} = \{Z_i \text{ má jakékoli spojité rozdělení}\}$

Testovaný parametr: Medián m_Z rozdílu $Z_i = X_i - Y_i$.

Hypotéza a alternativa:

$$H_0 : m_Z = 0, \quad H_1 : m_Z \neq 0.$$

Poznámka.

1. Medián Z_i obecně nelze vyjádřit pomocí mediánů X_i a Y_i . Obecně to tedy není test shody mediánů X_i a Y_i .
2. H_0 platí právě když $P[X_i \leq Y_i] = P[X_i \geq Y_i] = 1/2$, tj. X_i je s poloviční pravděpodobností větší než Y_i a s poloviční pravděpodobností menší než Y_i . Tj. jako test hypotézy nulového efektu (5.4) bude test konzistentní, pokud se vliv ošetření promítne do rozdělení Y_i tak, že $P[X_i \geq Y_i] \neq P[X_i \leq Y_i]$.
3. Pokud bychom zobecnili nulovou hypotézu a alternativu na

$$H_0 : m_Z = m_0, \quad H_1 : m_Z \neq m_0,$$

tak vlastně testujeme, že $P[X_i \leq Y_i + m_0] = P[X_i \geq Y_i + m_0] = 1/2$.

4. Má-li navíc Z_i konečnou střední hodnotu a hustotu **symetrickou** kolem 0, pak musí platit $E Z_i = E X_i - E Y_i = 0$. Za těchto dodatečných předpokladů je H_0 ekvivalentní hypotéze o rovnosti středních hodnot X_i a Y_i .

Testová statistika:

$$B_n = \sum_{i=1}^n \mathbb{1}\{Z_i > 0\}, \quad (\text{tj. počet párů, kde } X_i > Y_i).$$

Přesné rozdělení testové statistiky za H_0 :

$$B_n \sim \text{Bi}(n, \frac{1}{2})$$

Kritický obor (přesný test): Viz jednovýběrový znaménkový test.

Asymptotické rozdělení testové statistiky za H_0 :

$$\frac{B_n - \frac{n}{2}}{\sqrt{\frac{n}{4}}} \stackrel{\text{as.}}{\sim} N(0, 1)$$

Kritický obor (asymptotický test):

$$H_0 \text{ zamítneme} \Leftrightarrow \left| \frac{B_n - \frac{n}{2}}{\sqrt{\frac{n}{4}}} \right| \geq u_{1-\alpha/2}.$$

Poznámka. Výhodou párového znaménkového testu je, že nevyžaduje vyčíslení rozdílu mezi X_i a Y_i . Stačí informace o tom, že X_i je „lepší“ než Y_i , resp. X_i je „horší“ než Y_i . Tento test je vhodný pro aplikace, v nichž může být určení konkrétních hodnot X_i a Y_i problematické.

5.9. PÁROVÝ WILCOXONŮV TEST

Párový Wilcoxonův test* porovnává střed symetrie δ_Z rozdělení Z_i s předem danou konstantou.

Model: $\mathcal{F} = \{Z_i \text{ má spojité rozdělení a s hustotou } f \text{ splňující } \exists \delta \in \mathbb{R} : f(\delta - x) = f(\delta + x) \quad \forall x \in \mathbb{R}\}$

Poznámka. Předpoklad o symetrické hustotě se týká rozdílů Z_i , nikoli původních pozorování X_i a Y_i . Pokud existují příslušné střední hodnoty, pak předpoklad dává $\delta_Z = E Z_i = E X_i - E Y_i$.

Testované parametry: Střed symetrie δ_Z

Hypotéza a alternativa:

$$H_0 : \delta_Z = \delta_0, \quad H_1 : \delta_Z \neq \delta_0,$$

kde δ_0 je předem daná konstanta (obvykle $\delta_0 = 0$).

Testová statistika:

$$W_n = \sum_{i \in I} R_i,$$

kde $I \subset \{1, \dots, n\}$ je množina všech indexů takových, že $Z_i^* \stackrel{\text{df}}{=} X_i - Y_i - \delta_0$ má kladné znaménko pro $i \in I$, a R_i je pořadí náhodné veličiny $|Z_i^*|$ mezi všemi $|Z_1^*|, \dots, |Z_n^*|$.

Vlastnosti testové statistiky a kritický obor: viz jednovýběrový Wilcoxonův test.

Poznámka.

1. Vzhledem k poznámce výše se párový Wilcoxonův test může interpretovat jako test středních hodnot. K testování hypotézy rovnosti středních hodnot je zpravidla vhodnější párový t -test, protože nevyžaduje symetrii hustoty.
2. Pro $\delta_0 = 0$ můžeme použít test na testování hypotézy nulového efektu (5.4). V tomto případě je navíc přirozené za nulové hypotézy předpokládat, že sdružené rozdělení náhodného vektoru $(X_i, Y_i)^\top$ je stejné jako rozdělení $(Y_i, X_i)^\top$. Odtud pak plyne, že za nulové hypotézy je rozdělení náhodné veličiny $Z_i = X_i - Y_i$ **symetrické** kolem nuly, tj. test bude dodržovat předepsanou hladinu. Je však nutné si uvědomit, že test bude konzistentní pouze proti alternativám, pro které je pseudo-medián Z_i (tj. medián $\frac{Z_1 + Z_2}{2}$) nenulový. Tedy test je konzistentní vůči alternativám, pro které

$$P[Z_1 + Z_2 \leq 0] \neq P[Z_1 + Z_2 \geq 0].$$

* Angl. *paired Wilcoxon test, Wilcoxon signed rank test*

Přípravné příklady ke zkoušce.

Vaše řešení „praktických úloh“ by mělo obsahovat matematický model, hypotézu, testovou statistiku a její přesné (či asymptotické) rozdělení za nulové hypotézy. Dále pak kritický obor a vzorec pro výpočet p -hodnoty. Mělo by být také řečeno, zda je daný test přesný nebo asymptotický.

1. Uvádí se, že rozdělení IQ v běžné populaci má směrodatnou odchylku 15. Podařilo se nám získat údaje o IQ u 158 náhodně vybraných členů strany *Autoritářsky Nacionalistického Obrození*. Navrhněte test (tj. definujte vhodný model pro data, nulovou a alternativní hypotézu, testovou statistiku, kritický obor a vzorec pro p -hodnotu) jehož cílem by bylo prokázat, že členové této strany jsou (co se týká IQ) homogennější skupina, než je běžná populace.
2. Podařilo se nám získat údaje o hrubé mzdě 300 náhodně vybraných absolventek MFF UK. Navrhněte test, kterým byste chtěli prokázat, že alespoň 75 % absolventek MFF UK má hrubou mzdu vyšší než 40 000 Kč (pro srovnání mediánová hrubá mzda žen v ČR v roce 2019 byla cca 29 000 Kč).
3. Máte k dispozici údaje o mzdách 500 zaměstnanců v jedné pojišťovně. U každého zaměstnance máte dva údaje: nástupní mzdu a pak mzdu po dvou letech působení ve firmě. Navrhněte testy, kterými bychom chtěli prokázat, že během prvních dvou let:
 - (a) očekávané zvýšení mzdy je více než 15 000 Kč;
 - (b) s pravděpodobností (alespoň) 90 % mzda zaměstnance vzroste o méně než 10 000 Kč.

Předpokládejte, že by se Vám pomocí vhodných testů podařilo prokázat (a) i (b). Nebylo by to v rozporu?
4. Máme údaje o výšce 30 náhodně vybraných studentek MFF UK. Uvádí se, že průměrná výška dospělých žen v ČR je 168 cm. Rádi bychom prokázali, že studentky MFF UK jsou v nějakém smyslu spíše vyšší než je obvyklé v běžné populaci. Navrhněte vhodný test a vysvětlete, co bychom prokázali, pokud bychom zamítli nulovou hypotézu.
5. Následující tabulka zachycuje počty bodů z testu z anglického jazyka, které získalo 10 zaměstnanců před jazykovým kurzem a po jazykovém kurzu.

Zaměstnanec	1	2	3	4	5	6	7	8	9	10
před kurzem	37	41	36	48	42	36	42	44	40	34
po kurzu	38	43	43	47	52	44	41	42	42	39

Jak byste otestovali, zda můžeme tvrdit, že jazykový kurz zlepšil úroveň jazykových schopností zaměstnanců? *Nejde o samotné provedení testu, ale o jeho podrobné vysvětlení (tj. definujte vhodný model pro data, nulovou a alternativní hypotézu, testovou statistiku a kritický obor).*

6. Chceme si ověřit, zda daný tabulkový procesor má dobrý generátor náhodných čísel z rovnoměrného rozdělení $R(0, 1)$. Za tímto účelem si v procesoru vygenerujeme 1 000 realizací náhodných čísel. Jak byste otestovali, zda generátor náhodných čísel je opravdu dobrý?
7. Rozmyslete si, proč nedává dobrý smysl uvažovat párový Kolmogorovův-Smirnovův test.

6. DVOUVÝBĚROVÉ PROBLÉMY PRO KVANTITATIVNÍ DATA

Mějme dva *nezávislé* náhodné výběry: nechť X_1, \dots, X_n je náhodný výběr s distribuční funkcí F_X a Y_1, \dots, Y_m je náhodný výběr s distribuční funkcí F_Y . Model \mathcal{F} specifikuje množinu uvažovaných distribučních funkcí F_X a F_Y . Máme daný parametr $\theta = t(F)$, jehož hodnotu chceme pro oba výběry porovnat. Označme si $\theta_X = t(F_X)$ a $\theta_Y = t(F_Y)$. Obvykle chceme testovat hypotézu $H_0 : \theta_X = \theta_Y$ proti alternativě $H_1 : \theta_X \neq \theta_Y$, případně sestojit intervalový odhad pro rozdíl $\theta_X - \theta_Y$.

Dvouvýběrový problém lze zformulovat i jiným způsobem. Mějme náhodný výběr z dvourozměrného rozdělení

$$\begin{pmatrix} Z_1 \\ I_1 \end{pmatrix}, \dots, \begin{pmatrix} Z_N \\ I_N \end{pmatrix},$$

kde Z_j jsou hodnoty nezávislých stejně rozdělených měření a I_j má alternativní rozdělení. Indikátor I_j určuje, do které z porovnávaných skupin j -té pozorování patří (jestliže $I_j = 0$, pak do první skupiny, jinak do druhé). Přeznačíme-li si měření Z_j na X_i anebo Y_i podle toho, do jaké skupiny dané pozorování patří

$$(X_1, \dots, X_n) \stackrel{\text{df}}{=} (Z_j : I_j = 0) \quad \text{a} \quad (Y_1, \dots, Y_m) \stackrel{\text{df}}{=} (Z_j : I_j = 1),$$

získáme dva nezávislé výběry podle první formulace problému. Chceme porovnat podmíněné rozdělení Z_j v obou skupinách, tj. zajímají nás podmíněné distribuční funkce $F_X(x) = P[Z_j \leq x | I_j = 0]$ a $F_Y(x) = P[Z_j \leq x | I_j = 1]$. Případně jejich parametry $\theta_X = t(F_X)$ a $\theta_Y = t(F_Y)$. Tato druhá formulace dvouvýběrového problému je totožná s první, až na to, že rozsahy výběrů n a m nejsou konstanty, ale náhodné veličiny s binomickým rozdělením ($n = \sum_{j=1}^N (1 - I_j) \sim \text{Bi}(N, 1 - p_I)$, kde $p_I = P(I_j = 1)$). **Analýzu však provádíme stejně, jako by rozsahy výběrů byly pevné.**

Data podle první formulace získáme tak, že si předem stanovíme, kolik měření z každé skupiny budeme mít, a pak napozorujeme požadovaný počet veličin pro každou skupinu zvlášť. Data podle druhé formulace vzniknou, pokud stanovíme celkový počet pozorování $N = n + m$, učiníme N pozorování a u každého pozorování teprve dodatečně určíme, do které skupiny patří.

Obě formulace se trochu liší v pojetí asymptotických výsledků. U druhé formulace stačí vzít $N \rightarrow \infty$. U první formulace potřebujeme $n \rightarrow \infty$ a $m \rightarrow \infty$, ale navíc ještě musíme předpokládat, že rozsahy obou výběrů konvergují do nekonečna stejně rychle, tj. $n/m \rightarrow q$, kde $0 < q < \infty$.

Všechny metody uváděné v této kapitole se hodí pro obě formulace dvouvýběrového problému.

6.1. DVOUVÝBĚROVÝ KOLMOGOROVŮV-SMIRNOVŮV TEST

Dvouvýběrový Kolmogorovův-Smirnovův test* je rozšířením jednovýběrového testu stejného názvu. Je to neparametrický test, funguje pro jakákoli dvě spojitá rozdělení.

Model: $\mathcal{F} = \{\text{všechna spojitá rozdělení}\}$

Testované parametry: celé distribuční funkce F_X a F_Y

Hypotéza a alternativa:

$$H_0 : F_X(x) = F_Y(x) \quad \forall x \in \mathbb{R}, \quad H_1 : \exists x \in \mathbb{R} : F_X(x) \neq F_Y(x). \quad (6.1)$$

Testujeme, zdali oba výběry pocházejí z téhož rozdělení. Tuto hypotézu budeme dále nazývat **hypotézou nulového rozdílu**.

Testová statistika:

$$K_{n,m} = \sup_{x \in \mathbb{R}} |\widehat{F}_X(x) - \widehat{F}_Y(x)|,$$

kde \widehat{F}_X je empirická distribuční funkce náhodného výběru X_1, \dots, X_n a \widehat{F}_Y je empirická distribuční funkce náhodného výběru Y_1, \dots, Y_m .

Tvrzení 6.1 Nechť X_1, \dots, X_n a Y_1, \dots, Y_m jsou nezávislé náhodné výběry ze **spojitého** rozdělení s distribuční funkcí F_0 . Potom

$$\sqrt{\frac{nm}{n+m}} K_{n,m} \xrightarrow{d} Z, \text{ pro } m, n \rightarrow \infty,$$

kde náhodná veličina Z má distribuční funkci G danou předpisem (5.1).

Poznámka.

- Hypotézu zamítneme, pokud se empirické distribuční funkce obou výběrů od sebe příliš liší, tj. pokud je testová statistika velká.
- Tvrzení 6.1 implikuje, že za platnosti hypotézy konverguje $\sqrt{\frac{nm}{n+m}} K_{n,m}$ v distribuci k náhodné veličině s distribuční funkcí $G(y)$, která je stejná jako u jednovýběrového Kolmogorovova-Smirnovova testu (viz tvrzení 5.2). Důležité je, že tato distribuční funkce nezávisí na **neznámé** skutečné (společné) distribuční funkci F_0 (tj. testová statistika je pivotální). To nám umožňuje určit kritickou hodnotu pro zamítání H_0 .

Kritický obor:

$$H_0 \text{ zamítneme} \Leftrightarrow \sqrt{\frac{nm}{n+m}} K_{n,m} \geq k_{1-\alpha}, \quad (6.2)$$

kde $k_{1-\alpha} = G^{-1}(1 - \alpha)$ je $(1 - \alpha)$ -kvantil rozdělení s distribuční funkcí G , viz (5.1).

Podle tvrzení 6.1 má tento test asymptotickou hladinu α .

Poznámka.

* Angl. *two-sample Kolmogorov-Smirnov test*

- Je možné spočítat i přesnou kritickou hodnotu dvouvýběrového Kolmogorovova-Smirnovova testu pro spojitá rozdělení s malými rozsahy výběru n, m .
- Všimněme si, že za alternativy pro $m, n \rightarrow \infty$,

$$K_{n,m} \xrightarrow{P} \sup_{x \in \mathbb{R}} |F_X(x) - F_Y(x)| > 0.$$

Odtud se dá odvodit, že test je konsistentní (proti jakékoliv alternativě). Test tedy reaguje na jakýkoli rozdíl v rozděleních obou skupin. Další výhodou testu je absence omezujících předpokladů. Nevýhodou tohoto testu je, že má malou sílu proti specifickým druhům porušení H_0 . Zajímá-li nás (nebo očekáváme-li) pouze určitý typ porušení H_0 (třeba rozdíl ve střední hodnotě), je lepší použít test, který je zaměřen přímo na určitý parametr.

- Za povšimnutí stojí, že testová statistika se nezmění, pokud všechna pozorování nejdříve ztransformujeme pomocí nějaké prosté funkce g . Dá se ukázat, že Kolmogorovův-Smirnovův dvouvýběrový test se dá přepsat jako *pořadový test*.

PORUŠENÍ PŘEDPOKLADŮ

Pokud výběry za nulové hypotézy pochází z diskrétního rozdělení (tj. F_0 z tvrzení 6.1 není spojitá), tak test s kritickým oborem (6.2) bude konzervativní.

Podobně pokud „diskrétnost“ vznikne v důsledku zaokrouhlování. V tomto případě je však zapotřebí předpokládat, že způsob zaokrouhlování je pro oba dva výběry stejný.

6.2. DVOUVÝBĚROVÝ t -TEST BEZ PŘEDPOKLADU SHODY ROZPTYLŮ

Dvouvýběrový t -test* porovnává **střední hodnoty** obou výběrů. Provedení tohoto testu se liší podle toho, zda předpokládáme (viz kapitola 6.3), či nepředpokládáme shodu rozptylů.

Model:

$$\mathcal{F} = \{F_X \in \mathcal{L}_+^2, F_Y \in \mathcal{L}_+^2\}.$$

Testované parametry: Střední hodnoty $\mu_X = E X_i$ a $\mu_Y = E Y_j$.

Hypotéza a alternativa:

$$H_0 : \mu_X = \mu_Y + \delta_0, \quad H_1 : \mu_X \neq \mu_Y + \delta_0. \quad (6.3)$$

Testujeme, zdali se střední hodnoty obou výběrů liší o δ_0 (obvykle se klade $\delta_0 = 0$).

Testová statistika:

$$\tilde{T}_{n,m} = \frac{\bar{X}_n - \bar{Y}_m - \delta_0}{\sqrt{S_X^2/n + S_Y^2/m}},$$

* Angl. *two-sample t-test*

kde \bar{X}_n, \bar{Y}_m jsou aritmetické průměry obou výběrů a S_X^2, S_Y^2 jsou výběrové rozptyly.

Poznámka. Testová statistika $\tilde{T}_{n,m}$ se dá pamatovat následovně. Všimněme si, že

$$\text{var}(\bar{X}_n - \bar{Y}_m) = \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}.$$

Přirozeným (a dokonce nestranným) odhadem tohoto rozptylu pak je $S_X^2/n + S_Y^2/m$.

Věta 6.2 Necht' X_1, \dots, X_n a Y_1, \dots, Y_m jsou nezávislé náhodné výběry z rozdělení se středními hodnotami μ_X a μ_Y a konečnými rozptyly. Pak

$$\frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sqrt{S_X^2/n + S_Y^2/m}} \xrightarrow{d} N(0, 1) \text{ pro } m, n \rightarrow \infty, \frac{n}{m} \rightarrow q \in (0, \infty).$$

Důkaz. Přepišme

$$\frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sqrt{S_X^2/n + S_Y^2/m}} = \frac{\sqrt{m}(\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y))}{\sqrt{S_X^2 \frac{m}{n} + S_Y^2}}.$$

Z konzistence výběrového rozptylu $S_X^2 \xrightarrow{P} \sigma_X^2, S_Y^2 \xrightarrow{P} \sigma_Y^2$ a tudíž díky větě o spojité transformaci (tvrzení 1.2) $\sqrt{S_X^2 \frac{m}{n} + S_Y^2} \xrightarrow{P} \sqrt{\sigma_X^2/q + \sigma_Y^2}$. Tedy s využitím Cramérový-Sluckého věty (tvrzení 1.3) stačí ukázat, že

$$\sqrt{m}(\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)) \xrightarrow{d} N(0, \sigma_X^2/q + \sigma_Y^2). \quad (6.4)$$

Z centrální limitní věty $\sqrt{n}(\bar{X}_n - \mu_X) \xrightarrow{d} N(0, \sigma_X^2)$ a tudíž

$$\sqrt{m}(\bar{X}_n - \mu_X) = \sqrt{\frac{m}{n}} \sqrt{n}(\bar{X}_n - \mu_X) \xrightarrow{d} N(0, \sigma_X^2/q), \quad (6.5)$$

neboť z předpokladů věty $\sqrt{\frac{m}{n}} \rightarrow \frac{1}{\sqrt{q}}$. Dále díky centrální limitní větě

$$\sqrt{m}(\bar{Y}_m - \mu_Y) \xrightarrow{d} N(0, \sigma_Y^2). \quad (6.6)$$

Nyní s využitím (6.5), (6.6) a **nezávislosti** \bar{X}_n a \bar{Y}_m

$$\sqrt{m} \begin{pmatrix} \bar{X}_n - \mu_X \\ \bar{Y}_m - \mu_Y \end{pmatrix} \xrightarrow{d} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_X^2/q & 0 \\ 0 & \sigma_Y^2 \end{pmatrix} \right).$$

Tudíž také pro všechny $c \in \mathbb{R}^2$

$$c^T \sqrt{m} \begin{pmatrix} \bar{X}_n - \mu_X \\ \bar{Y}_m - \mu_Y \end{pmatrix} \xrightarrow{d} N(0, c^T \Sigma c), \quad \text{kde } \Sigma = \begin{pmatrix} \sigma_X^2/q & 0 \\ 0 & \sigma_Y^2 \end{pmatrix}.$$

Konvergence v (6.4) nyní plyne z výše uvedené konvergence pro $c = (1, -1)^T$, neboť v tomto případě

$$c^T \Sigma c = \sigma_X^2/q + \sigma_Y^2.$$

□

Poznámka.

- Hypotézu budeme zamítat, pokud se výběrové průměry obou skupin od sebe příliš liší, tj. pokud je testová statistika buď moc velká nebo moc malá.
- Věta 6.2 implikuje, že za platnosti modelu \mathcal{F} a hypotézy H_0 má $\tilde{T}_{n,m}$ asymptoticky rozdělení $N(0, 1)$.

Kritický obor:

$$H_0 \text{ zamítneme} \Leftrightarrow \left| \tilde{T}_{n,m} \right| \geq u_{1-\alpha/2},$$

kde $u_{1-\alpha/2}$ je $(1 - \alpha/2)$ -tý kvantil normovaného normálního rozdělení.

P-hodnota: $p = 2(1 - \Phi(|t|))$, kde t je pozorovaná hodnota testové statistiky $\tilde{T}_{n,m}$ a Φ je distribuční funkce rozdělení $N(0, 1)$.

Interval spolehlivosti pro $\mu_X - \mu_Y$: Z věty 6.2 lze odvodit asymptotický interval spolehlivosti pro rozdíl středních hodnot obou výběrů. Pro $n, m \rightarrow \infty$ dostáváme

$$P \left[\bar{X}_n - \bar{Y}_m - u_{1-\alpha/2} \sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}} < \mu_X - \mu_Y < \bar{X}_n - \bar{Y}_m + u_{1-\alpha/2} \sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}} \right] \rightarrow 1 - \alpha.$$

Ukazuje se, že i kdybychom přidali předpoklad normality našich pozorování, tj. $X_i \sim N(\mu_X, \sigma_X^2)$ a $Y_i \sim N(\mu_Y, \sigma_Y^2)$, tak rozdělení testové statistiky $\tilde{T}_{n,m}$ je za hypotézy stále pouze asymptoticky pivotální. Přesné rozdělení statistiky $\tilde{T}_{n,m}$ i za předpokladu normality bude totiž záviset na poměru σ_X^2/σ_Y^2 . Proto se v praxi spokojujeme s asymptotickým testem.*

Poznámka. Existují i lepší aproximace kritických hodnot pro tento test založené na t -rozdělení s počtem stupňů volnosti, který závisí na počtu pozorování v obou skupinách a výběrových rozptylech. Takových aproximací je několik[†]. Jedna z variant této aproximace, tzv. *Welchův test*[‡], je implementována ve výpočetním prostředí R jako standardní metoda testování rovnosti středních hodnot dvou výběrů (provádí jej funkce `t.test`). V této variantě se používají jako kritické hodnoty kvantily t -rozdělení se stupni volnosti f danými vzorcem

$$f = \frac{\left(\frac{S_X^2}{n} + \frac{S_Y^2}{m} \right)^2}{\frac{(S_X^2)^2}{n^2(n-1)} + \frac{(S_Y^2)^2}{m^2(m-1)}}.$$

Tento vzorec byl odvozen za předpokladu normality dat a na základě aproximace rozdělení náhodné veličiny $\frac{S_X^2}{n} + \frac{S_Y^2}{m}$ z čitatele testové statistiky pomocí násobku χ^2 -rozdělení s „vhodným“ počtem stupňů volnosti (detaily lze nalézt ve [Welch, 1938](#)).

Welchův test lze tedy chápat jako variantu dvouvýběrového t -testu (bez předpokladu shody rozptylů) s vylepšenými kritickými hodnotami. P-hodnotu Welchova t -testu pro oboustrannou alternativu (6.3) pak spočteme jako $p = 2(1 - G_f(|t|))$, kde t

* Problém provedení přesného testu je znám jako tzv. Behrensův-Fisherův problém. [†] Lze je nalézt např. v kapitole 8.1. knihy [Anděl \(1998\)](#). [‡] Angl. *Welch test*

je pozorovaná hodnota testové statistiky $\tilde{T}_{n,m}$ a G_f je distribuční funkce t -rozdělení o f stupních volnosti.

6.3. DVOUVÝBĚROVÝ t -TEST ZA PŘEDPOKLADU SHODY ROZPTYLŮ

Podobně jako v případě jednovýběrového t -testu (viz kapitola 5.3) bude tento test přesný za předpokladu normality pozorování a asymptotický bez tohoto předpokladu.

Model:

$$\mathcal{F}_n = \{F_X = N(\mu_X, \sigma^2), F_Y = N(\mu_Y, \sigma^2), \mu_X, \mu_Y \in \mathbb{R}, \sigma^2 > 0\}$$

nebo

$$\mathcal{F}_{as} = \{F_X \in \mathcal{L}_+^2, F_Y \in \mathcal{L}_+^2, \text{přičemž } \text{var}(X_i) = \text{var}(Y_i) \stackrel{\text{df}}{=} \sigma^2\}.$$

V modelu \mathcal{F}_n mají tedy oba výběry normální rozdělení s totožným rozptylem σ^2 , tj. mohou se lišit pouze střední hodnotou. V modelu \mathcal{F}_{as} požadujeme pouze, a aby byly totožné rozptyly (tj. rozdělení mohou být různá). Všimněme si, že $\mathcal{F}_n \subset \mathcal{F}_{as}$.

Testované parametry: Střední hodnoty $\mu_X = E X_i$ a $\mu_Y = E Y_j$.

Hypotéza a alternativa:

$$H_0 : \mu_X = \mu_Y + \delta_0, \quad H_1 : \mu_X \neq \mu_Y + \delta_0.$$

Testujeme, zdali se střední hodnoty obou výběrů liší o δ_0 (obvykle se klade $\delta_0 = 0$).

Testová statistika:

$$T_{n,m} = \frac{\bar{X}_n - \bar{Y}_m - \delta_0}{\sqrt{S_{n,m}^2 \left(\frac{1}{n} + \frac{1}{m}\right)}} = \sqrt{\frac{nm}{n+m}} \frac{\bar{X}_n - \bar{Y}_m - \delta_0}{S_{n,m}},$$

kde \bar{X}_n a \bar{Y}_m jsou aritmetické průměry obou výběrů a

$$S_{n,m}^2 \stackrel{\text{df}}{=} \frac{1}{n+m-2} \left[\sum_{i=1}^n (X_i - \bar{X}_n)^2 + \sum_{j=1}^m (Y_j - \bar{Y}_m)^2 \right] = \frac{n-1}{n+m-2} S_X^2 + \frac{m-1}{n+m-2} S_Y^2$$

je nestranný odhad společného rozptylu σ^2 spočítaný z obou výběrů (vážený průměr výběrových rozptylů S_X^2 a S_Y^2).

Poznámka. Testovou statistiku $T_{n,m}$ si lze pamatovat následovně. Všimněme si, že

$$\text{var}(\bar{X}_n - \bar{Y}_m) = \frac{\sigma^2}{n} + \frac{\sigma^2}{m} = \sigma^2 \left(\frac{1}{n} + \frac{1}{m}\right).$$

Jelikož $S_{n,m}^2$ je (nestranný) odhad σ^2 , tak $S_{n,m}^2 \left(\frac{1}{n} + \frac{1}{m}\right)$ je přirozený (a dokonce nestranný) odhad $\sigma^2 \left(\frac{1}{n} + \frac{1}{m}\right)$.

Věta 6.3 Necht' X_1, \dots, X_n a Y_1, \dots, Y_m jsou nezávislé náhodné výběry z rozdělení se středními hodnotami μ_X a μ_Y a konečnými rozptyly $\sigma_X^2 = \text{var}(X_i)$ a $\sigma_Y^2 = \text{var}(Y_j)$. Pak

$$\frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sqrt{S_{n,m}^2 \left(\frac{1}{n} + \frac{1}{m}\right)}} \xrightarrow{d} N(0, \sigma_*^2), \text{ pro } m, n \rightarrow \infty, \frac{n}{n+m} \rightarrow \lambda \in (0, 1),$$

kde

$$\sigma_*^2 = \frac{(1 - \lambda)\sigma_X^2 + \lambda\sigma_Y^2}{\lambda\sigma_X^2 + (1 - \lambda)\sigma_Y^2}.$$

Důkaz. Důkaz je analogický důkazu věty 6.2. Nejdříve si přepíšme

$$\frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sqrt{S_{n,m}^2 \left(\frac{1}{n} + \frac{1}{m}\right)}} = \frac{\sqrt{m} (\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y))}{\sqrt{S_{n,m}^2 \frac{m+n}{n}}}.$$

Nyní stejně jako ve větě 6.2 ukážeme, že

$$\sqrt{m} (\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)) \xrightarrow{d} N\left(0, \frac{1-\lambda}{\lambda} \sigma_X^2 + \sigma_Y^2\right),$$

kde jsme pouze využili toho, že $\frac{1}{q} = \frac{1-\lambda}{\lambda}$. Dále pak ukážeme, že

$$\sqrt{S_{n,m}^2 \frac{m+n}{n}} \xrightarrow[n \rightarrow \infty]{P} \sqrt{\frac{\lambda\sigma_X^2 + (1-\lambda)\sigma_Y^2}{\lambda}}.$$

□

Všimněme si, že pokud platí předpoklad shody rozptylů, tj. $\sigma_X^2 = \sigma_Y^2$ (tj. platí model \mathcal{F}_{as}), pak $\sigma_*^2 = 1$ a tedy

$$\frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sqrt{S_{n,m}^2 \left(\frac{1}{n} + \frac{1}{m}\right)}} \xrightarrow{d} N(0, 1), \text{ pro } m, n \rightarrow \infty, \frac{n}{n+m} \rightarrow \lambda \in (0, 1).$$

V případě, že mohu navíc přidat předpoklad *normality* (tj. platí model \mathcal{F}_n), pak můžeme odvodit přesné rozdělení náhodné veličiny $\frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sqrt{S_{n,m}^2 \left(\frac{1}{n} + \frac{1}{m}\right)}}$.

Věta 6.4 Necht' X_1, \dots, X_n a Y_1, \dots, Y_m jsou nezávislé náhodné výběry z **normálních** rozdělení se středními hodnotami μ_X a μ_Y a se **shodným rozptylem** σ^2 . Pak

$$\frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sqrt{S_{n,m}^2 \left(\frac{1}{n} + \frac{1}{m}\right)}} \sim t_{n+m-2}.$$

Důkaz. Přepišme

$$\frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sqrt{S_{n,m}^2(\frac{1}{n} + \frac{1}{m})}} = \frac{U}{\sqrt{Z/(n+m-2)}},$$

kde

$$U = \frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sigma\sqrt{\frac{1}{n} + \frac{1}{m}}} \quad \text{a} \quad Z = \frac{(n+m-2)S_{n,m}^2}{\sigma^2}.$$

K dokončení důkazu stačí ukázat, že $U \sim N(0, 1)$, $Z \sim \chi_{n+m-2}^2$ a U je nezávislé se Z .

1. $U \sim N(0, 1)$. K tomu si stačí uvědomit, že díky nezávislosti náhodných výběrů jsou také \bar{X}_n a \bar{Y}_m nezávislé a platí

$$\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y) \sim N(0, \frac{\sigma^2}{n} + \frac{\sigma^2}{m}).$$

Tedy

$$U = \frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sigma\sqrt{\frac{1}{n} + \frac{1}{m}}} \sim N(0, 1).$$

2. $Z \sim \chi_{n+m-2}^2$. Náhodnou veličinu Z si můžeme napsat pomocí S_X^2 a S_Y^2 jako

$$Z = \frac{(n+m-2)S_{n,m}^2}{\sigma^2} = \frac{(n-1)S_X^2}{\sigma^2} + \frac{(m-1)S_Y^2}{\sigma^2}.$$

Nyní díky větě 2.8(i) víme, že $\frac{(n-1)S_X^2}{\sigma^2} \sim \chi_{n-1}^2$ a $\frac{(m-1)S_Y^2}{\sigma^2} \sim \chi_{m-1}^2$.

Dále z nezávislosti S_X^2 a S_Y^2 dostáváme, že rozdělení Z je stejné jako rozdělení součtu dvou **nezávislých** veličin s rozděleními χ_{n-1}^2 a χ_{m-1}^2 . Nyní z definice χ^2 -rozdělení (jako rozdělení součtu kvadrátů nezávislých $N(0, 1)$ veličin) plyne, že $Z \sim \chi_{n+m-2}^2$.

3. Nezávislost U a Z . Díky nezávislosti náhodných výběrů jsou náhodné vektory (\bar{X}_n, S_X^2) a (\bar{Y}_m, S_Y^2) nezávislé. Dále z větě 2.8(ii) jsou náhodné veličiny \bar{X}_n a S_X^2 nezávislé a podobně také náhodné veličiny \bar{Y}_m a S_Y^2 jsou nezávislé. Tudíž také náhodné veličiny $\bar{X}_n - \bar{Y}_m$ a $S_{n,m}^2$ jsou nezávislé. Odtud již plyne nezávislost U a Z . \square

Poznámka.

- Zatímco z věty 6.3 plyne, že za platnosti modelu \mathcal{F}_{as} má $T_{n,m}$ asymptoticky rozdělení $N(0, 1)$, tak věta 6.4 nám říká, že v menším modelu \mathcal{F}_n má $T_{n,m}$ za H_0 přesně rozdělení t_{n+m-2} .
- Hypotézu budeme zamítat, pokud se výběrové průměry obou skupin od sebe příliš liší, tj. pokud je testová statistika buď moc velká nebo moc malá.

Kritický obor:

$$H_0 \text{ zamítáme} \Leftrightarrow |T_{n,m}| \geq t_{n+m-2}(1 - \alpha/2),$$

kde $t_{n+m-2}(1 - \alpha/2)$ je $(1 - \alpha/2)$ -tý kvantil t -rozdělení s $n + m - 2$ stupni volnosti.

Podobně jako u jednovýběrového t -testu je výše uvedený test přesný v modelu \mathcal{F}_n a asymptotický v modelu \mathcal{F}_{as} . Totéž platí o níže uvedené p -hodnotě a intervalu spolehlivosti.

P -hodnota: $p = 2(1 - F(|t|))$, kde t je pozorovaná hodnota testové statistiky $T_{n,m}$ a F je distribuční funkce rozdělení t_{n+m-2} .

Interval spolehlivosti pro $\mu_X - \mu_Y$: Z věty 6.4 (resp. věty 6.3) lze odvodit přesný (resp. asymptotický) interval spolehlivosti pro rozdíl středních hodnot obou výběrů. Dostaneme

$$P\left[\bar{X}_n - \bar{Y}_m - t_{n+m-2}(1 - \alpha/2) S_{n,m} \sqrt{\frac{1}{n} + \frac{1}{m}} < \mu_X - \mu_Y < \bar{X}_n - \bar{Y}_m + t_{n+m-2}(1 - \alpha/2) S_{n,m} \sqrt{\frac{1}{n} + \frac{1}{m}}\right] = 1 - \alpha.$$

Cvičení. Upravte kritický obor a vzorec pro výpočet p -hodnoty pro test hypotézy $H_0 : \mu_X \leq \mu_Y + \delta_0$ proti alternativě $H_1 : \mu_X > \mu_Y + \delta_0$.

PORUŠENÍ PŘEDPOKLADU SHODY ROZPTYLŮ

Dle věty 6.3

$$T_{n,m} \xrightarrow{d} N(0, \sigma_*^2), \quad \text{kde} \quad \sigma_*^2 = \frac{(1 - \lambda)\sigma_X^2 + \lambda\sigma_Y^2}{\lambda\sigma_X^2 + (1 - \lambda)\sigma_Y^2} \quad \text{a} \quad \frac{n}{n+m} \rightarrow \lambda \in (0, 1).$$

Tedy test obecně nedodrží předepsanou hladinu ani asymptoticky. Přičemž stojí za povšimnutí, že pokud například $\sigma_X^2 > \sigma_Y^2$ a zároveň $\lambda < \frac{1}{2}$ (tj. větší rozptyl je ve výběru s menším rozsahem), pak $\sigma_*^2 > 1$ a test je (asymptoticky) anti-konzervativní.

Všimněme si také, že pro $\lambda = \frac{1}{2}$ je $\sigma_*^2 = 1$. Tedy (pro přibližně) stejné rozsahy výběru stále platí

$$T_{n,m} \xrightarrow{d} N(0, 1)$$

a test asymptoticky dodržuje hladinu.

Navíc pokud dokonce oba výběry mají stejný rozsah, tj. $m = n$, pak

$$\sqrt{S_X^2/n + S_Y^2/m} = \sqrt{(S_X^2/2 + S_Y^2/2) \frac{2}{n}} = \sqrt{S_{n,m}^2 \left(\frac{1}{n} + \frac{1}{m}\right)}.$$

V tomto případě tedy vždy platí $T_{n,n} = \tilde{T}_{n,n}$, tj. testové statistiky dvouvýběrového t -testu jsou totožné, ať již předpokládáme stejné rozptyly či nikoliv.

***t*-TEST JAKO TEST HYPOTÉZY NULOVÉHO ROZDÍLU**

Pro $\delta_0 = 0$ můžeme na test nahlížet jako na test hypotézy nulového rozdílu (6.1). V tomto případě sice nemáme zaručenu normalitu, ale za nulové hypotézy jsou rozptyly shodné. Test tedy bude dodržovat hladinu asymptoticky.

Co se týká síly testu, tak test bude konsistentní vůči alternativám, pro které $\mu_X - \mu_Y \neq 0$. Pokud se však kromě změny střední hodnoty budou rozdělení F_X a F_Y lišit také rozptylem, tak vliv rozdílnosti těchto rozptylů nemáme pod kontrolou. Rozdílné rozptyly mohou zvyšovat či snižovat sílu testu. Navíc při zamítnutí nulové hypotézy (6.1) můžeme pouze tvrdit, že jsme prokázali rozdílnost rozdělení F_X a F_Y . Toto zamítnutí však nemůžeme přisuzovat pouze k rozdílu středních hodnot, protože k němu mohla přispět i rozdílnost rozptylů.

Cvičení. Podrobně dokažte větu 6.3.

Poznámka. Někdy se doporučuje před použitím dvouvýběrového *t*-testu otestovat shodu rozptylů obou výběrů, např. testem uvedeným v kap. 6.5 níže, nebo tzv. Leveneovým testem (neuvádíme). Pokud test zamítne rovnost rozptylů, použijeme Welchův test, jinak použijeme dvouvýběrový *t*-test. Od používání takového postupu spíše odrazujeme. Jedná se o tzv. *dvoufázový test*, kdy celkový výsledek testu závisí na třech různých vzájemně závislých testových statistikách. Není ničím zaručeno, že celková hladina takové testovací procedury je rovna požadované hodnotě α . Pokud si nejsme jisti shodou rozptylů nebo normalitou dat, provedeme raději rovnou Welchův test. Ani jeden z předpokladů dvouvýběrového *t*-testu pak není třeba nijak ověřovat.

6.4. DVOUVÝBĚROVÝ WILCOXONŮV TEST

Dvouvýběrový Wilcoxonův test* je neparametrický test založený na pořadích.

Model: $\mathcal{F} = \{ \exists \text{ rostoucí funkce } g \text{ a } \exists \delta \in \mathbb{R} :$

$$g(X_i) \sim \tilde{F}_X \text{ spojitá d.f., } g(Y_j) \sim \tilde{F}_Y, \tilde{F}_X(x) = \tilde{F}_Y(x - \delta) \forall x \in \mathbb{R} \}. \quad (6.7)$$

Testovaný parametr: Posunutí δ_{XY} .

Hypotéza a alternativa:

$$H_0 : \delta_{XY} = 0, \quad H_1 : \delta_{XY} \neq 0.$$

Pokud funkce $g(x) = x$, pak model \mathcal{F} nazýváme *model posunutí v poloze*.† Model \mathcal{F} tedy budeme nazývat *zobecněný model posunutí*.

Poznámka.

- Na rozdíl od jednovýběrového a párového Wilcoxonova testu **nevyžadujeme symetrii** žádné z hustot.

* Angl. *two-sample Wilcoxon test, Wilcoxon rank-sum test* † Angl. *location model*.

- Pokud platí model \mathcal{F} a hypotéza H_0 , rozdělení X a Y jsou **totožná**. Potom platí $m_X = m_Y$ a $E X = E Y$ (existují-li střední hodnoty). To jest, za platnosti modelu \mathcal{F} lze dvouvýběrový Wilcoxonův test chápat jako test rovnosti středních hodnot i mediánů. Většinou se však o dvouvýběrovém Wilcoxonově testu mluví jako o testu shody mediánů.

Testová statistika:

$$W_{n,m} = \sum_{i=1}^n R_i, \quad (6.8)$$

kde R_1, R_2, \dots, R_n jsou pořadí náhodných veličin X_i ve spojeném náhodném výběru $X_1, \dots, X_n, Y_1, \dots, Y_m$.

Poznámka. Testová statistika $W_{n,m}$ může nabývat hodnot z množiny $\{\frac{n(n+1)}{2}, \dots, mn + \frac{n(n+1)}{2}\}$. Statistiku spočítáme následujícím způsobem:

1. Vezmeme spojený výběr $(Z_1, \dots, Z_{n+m}) \stackrel{\text{df}}{=} (X_1, \dots, X_n, Y_1, \dots, Y_m)$.
2. Seřadíme všechny Z_j od nejmenší do největší; získáme uspořádaný výběr

$$Z_{(1)} < Z_{(2)} < \dots < Z_{(n+m)}.$$

3. Určíme pořadí R_i náhodné veličiny X_i mezi všemi $Z_{(1)}, \dots, Z_{(n+m)}$. Platí $X_i = Z_{(R_i)}$.
4. Sečteme pořadí R_i pro $i = 1, \dots, n$.

Nejsou-li n a m příliš velká, lze nalézt i *přesné rozdělení* testové statistiky $W_{n,m}$ za nulové hypotézy (numericky nebo v tabulkách). Toto přesné rozdělení se dá odvodit z toho, že za **nulové hypotézy** je jakékoliv uspořádání veličin Z_1, \dots, Z_{n+m} stejně pravděpodobné (viz věta 2.15) a tudíž

$$P(R_1 = r_1, \dots, R_n = r_n) = \frac{m!}{(n+m)!}$$

pro všechna $r_1, \dots, r_n \in \{1, \dots, n+m\}$ různá.

Pro velké hodnoty n a m se využívá následující tvrzení.

Tvrzení 6.5 Nechť X_1, \dots, X_n a Y_1, \dots, Y_m jsou nezávislé náhodné výběry, pro které platí model \mathcal{F} . Nechť dále platí **hypotéza** H_0 , pak

(i)

$$E_{H_0} W_{n,m} = \frac{n(n+m+1)}{2}, \quad \text{var}_{H_0}(W_{n,m}) = \frac{mn(n+m+1)}{12}.$$

(ii) Pokud $n, m \rightarrow \infty$,

$$\frac{W_{n,m} - E_{H_0} W_{n,m}}{\sqrt{\text{var}_{H_0}(W_{n,m})}} \xrightarrow{d} N(0, 1).$$

Důkaz. Část (i). Za platnosti hypotézy jsou rozdělení X_i a Y_j shodné, tedy spojený výběr $X_1, \dots, X_n, Y_1, \dots, Y_m$ je náhodný výběr o rozsahu $n+m$. S využitím věty 2.16 tedy máme, že pro $i \neq j$:

$$E R_i = \frac{n+m+1}{2}, \quad \text{var}(R_i) = \frac{(n+m)^2 - 1}{12}, \quad \text{cov}(R_i, R_j) = -\frac{n+m+1}{12}.$$

Tedy

$$E_{H_0} W_{n,m} = \sum_{i=1}^n E R_i = \frac{n(n+m+1)}{2}$$

a

$$\begin{aligned} \text{var}_{H_0}(W_{n,m}) &= \sum_{i=1}^n \text{var}(R_i) + \sum_{i=1}^n \sum_{j=1, j \neq i}^n \text{cov}(R_i, R_j) \\ &= n \frac{(n+m+1)(n+m-1)}{12} - n(n-1) \frac{n+m+1}{12} \\ &= \frac{n(n+m+1)}{12} [n+m-1 - (n-1)] = \frac{nm(n+m+1)}{12}. \end{aligned}$$

Část (ii). Nebudeme dokazovat. Potíž důkazu spočívá v tom, že pořadí R_1, \dots, R_n nejsou nezávislé náhodné veličiny. \square

Poznámka.

- Hypotézu budeme zamítat pro příliš malé nebo příliš velké hodnoty $W_{n,m}$.
- Předchozí tvrzení dává návod k nalezení kritických hodnot pro zamítání hypotézy, které zaručují asymptotickou hladinu α .

Kritický obor (asymptotický test):

$$H_0 \text{ zamítneme} \Leftrightarrow \frac{\left| W_{n,m} - \frac{n(m+n+1)}{2} \right|}{\sqrt{\frac{mn(m+n+1)}{12}}} \geq u_{1-\alpha/2}.$$

PORUŠENÍ PŘEDPOKLADŮ

Shody kvůli zaokrouhlování. Kvůli zaokrouhlování bývají v aplikacích v datech často shody. Testová statistika $W_{n,m}$ se pak spočte s využitím tzv. *průměrných pořadí*, viz. (2.6). Dá se ukázat, že za nulové hypotézy

$$\frac{W_{n,m} - \frac{n(m+n+1)}{2}}{\sqrt{\frac{mn(m+n+1 - \text{kor.})}{12}}} \xrightarrow{d} N(0, 1), \quad \text{pro } n, m \rightarrow \infty.$$

kde $kor.$ je korekce upravující rozptyl daná předpisem*

$$kor. = \frac{1}{(n+m)(n+m-1)} \sum_z (t_z^3 - t_z),$$

kde t_z značí počet, kolikrát se mezi hodnotami Z_1, \dots, Z_{n+m} vyskytla hodnota z . Suma \sum_z pak značí sčítání přes všechny rozdílné hodnoty množiny $\{Z_1, \dots, Z_{n+m}\}$.

Za povšimnutí stojí, že bez úpravy jmenovatele pomocí $kor.$ by byl test asymptoticky konzervativní.

Neplatí zobecněný model posunutí \mathcal{F} . Nejdříve si všimněme, že test (asymptoticky) dodržuje hladinu, pokud platí hypotéza nulového rozdílu, tj. $F_X = F_Y$. Neplatnost modelu posunutí má tedy vliv na interpretaci a sílu testu.

Co se týká **interpretace testu**, tak zamítnutí nulové hypotézy nám mimo model zobecněného posunutí pouze říká, že rozdělení F_X a F_Y nejsou totožná. Obecně však nelze tvrdit, že se liší mediány, resp. střední hodnoty těchto rozdělení.

Co se týká **síly**, tak ve výše popsáném zobecněném modelu posunutí je Wilcoxonův test konsistentní.

V praxi však nikdy nemáme jistotu, zda zobecněný model posunutí platí. Proto pro hlubší porozumění dvouvýběrovému Wilcoxonovu testu je vhodné využít níže uvedenou Mannovu-Whitneyho formulaci Wilcoxonova testu.

MANNOVA-WHITNEYHO FORMULACE WILCOXONOVA TESTU

Test ekvivalentní s Wilcoxonovým lze získat i následující úvahou. Uvažujme všechny dvojice (X_i, Y_j) pro $i \in \{1, \dots, n\}$ a $j \in \{1, \dots, m\}$ a spočtěme, kolik z nich splňuje podmínku $X_i < Y_j$:

$$W_{n,m}^* = \sum_{i=1}^n \sum_{j=1}^m \mathbb{1}\{X_i < Y_j\}.$$

Náhodná veličina $W_{n,m}^*$, tzv. *Mannova-Whitneyho statistika*, může nabývat hodnot z množiny $\{0, \dots, nm\}$.

Následující tvrzení ukazuje, že mezi dvouvýběrovou Wilcoxonovou statistikou $W_{n,m}$ a Mannovou-Whitneyho statistikou $W_{n,m}^*$ je deterministický lineární vztah. Můžeme tedy snadno spočítat momenty $W_{n,m}^*$.

Tvrzení 6.6

(i) $W_{n,m} + W_{n,m}^* = nm + \frac{n(n+1)}{2}$.

(ii) Pokud $\min(n, m) \rightarrow \infty$, pak $\frac{W_{n,m}^*}{nm} \xrightarrow{P} P[X_i < Y_j]$.

Důkaz. Část (i). Z definice pořadí

$$R_i = \sum_{j=1}^{n+m} \mathbb{1}\{Z_j \leq X_i\} = \sum_{j=1}^n \mathbb{1}\{X_j \leq X_i\} + \sum_{j=1}^m \mathbb{1}\{Y_j \leq X_i\}.$$

* Viz např. [Hollander et al. \(2013\)](#), str. 118.

Tedy

$$\begin{aligned}
 W_{n,m} + W_{n,m}^* &= \sum_{i=1}^n R_i + \sum_{i=1}^n \sum_{j=1}^m \mathbb{1}\{X_i < Y_j\} \\
 &= \sum_{i=1}^n \sum_{j=1}^n \mathbb{1}\{X_j \leq X_i\} + \sum_{i=1}^n \sum_{j=1}^m \mathbb{1}\{Y_j \leq X_i\} + \sum_{i=1}^n \sum_{j=1}^m \mathbb{1}\{X_i < Y_j\} \\
 &= \sum_{i=1}^n \sum_{j=1}^n \mathbb{1}\{X_{(j)} \leq X_{(i)}\} + \sum_{i=1}^n \sum_{j=1}^m \mathbb{1}\{Y_j \leq X_i \text{ nebo } Y_j > X_i\} \\
 &= \sum_{i=1}^n i + nm = \frac{n(n+1)}{2} + nm.
 \end{aligned}$$

Část (ii). Nebudeme dokazovat. Potíž důkazu vězí v tom, že indikátory $\mathbb{1}\{X_i < Y_j\}$ nejsou (pro $i \in \{1, \dots, n\}$, $j \in \{1, \dots, m\}$) nezávislé náhodné veličiny. \square

Rozeberme si důsledky tvrzení 6.6. Část (i) říká, že testy založené na dvouvýběrové Wilcoxonově statistice a Mannově-Whitneyho statistice jsou ekvivalentní. Části (ii) pak ukazuje, že $\frac{W_{n,m}^*}{nm}$ je konsistentním odhadem parametru $\theta_{XY} = P[X_i < Y_j]$. Pokud $F_X = F_Y$, lze snadno ukázat, že $\theta_{XY} = 1/2$. Parametr θ_{XY} však může nabývat hodnoty $1/2$ i pro dvě rozdělení, která nejsou totožná.

Tedy uvažujeme-li dvouvýběrový Wilcoxonův test jako test hypotézy nulového rozdílu (6.1), pak je tento test konsistentní pouze vůči alternativám, pro které je $\theta_{XY} \neq \frac{1}{2}$. Tuto nerovnost však obecně (tj. mimo model posunutí) nemůžeme interpretovat jako nerovnost středních hodnot nebo mediánů. Existují spojitá rozdělení F_X a F_Y taková, že mají rozdílné střední hodnoty (resp. mediány) a přitom $\theta_{XY} = \frac{1}{2}$. A na druhou stranu existují spojitá rozdělení F_X a F_Y taková, že mají stejné střední hodnoty (resp. mediány) a přitom $\theta_{XY} \neq \frac{1}{2}$.

Vzhledem k výše uvedenému by nás mohlo zajímat, zda bychom nemohli uvažovat Mannův-Whitneyho test jako test pro následující obecnou situaci.

Model: $\mathcal{F}^* = \{X \sim F_X \text{ spojitá d.f.}, Y \sim F_Y \text{ spojitá d.f.}\}$

Testovaný parametr: $\theta_{XY} = P[X < Y]$

Hypotéza a alternativa:

$$H_0^* : \theta_{XY} = \frac{1}{2}, \quad H_1^* : \theta_{XY} \neq \frac{1}{2}.$$

Problém však je, že v tomto případě nelze rozptýlit testové statistiky $W_{n,m}^*$ za hypotézy počítat podle tvrzení 6.5 (neboť za hypotézy již obecně nemáme stejně rozdělené náhodné veličiny). Kritické hodnoty spočítané pro Wilcoxonův test v modelu \mathcal{F} tedy v obecném modelu \mathcal{F}^* nefungují. Přičemž se ukazuje, že ignorování tohoto faktu může vést k testu, který je konzervativní nebo naopak anti-konzervativní.*

* Standardizaci testové statistiky $W_{n,m}^*$, aby test asymptoticky dodržoval hladinu i v obecném modelu \mathcal{F}^* , lze nalézt například v [Chung and Romano \(2016\)](#).

Tyto úvahy vedou k jednoznačnému závěru: *Chceme-li testovat rovnost středních hodnot bez dalších předpokladů na tvar rozdělení obou výběrů, použijeme dvouvýběrový t -test bez předpokladu shody rozptylů (Welchův test), nikoli Wilcoxonův test.*

Poznámka. Někdy se doporučuje před použitím dvouvýběrového t -testu na porovnání středních hodnot otestovat normalitu obou výběrů (populární je například tzv. Shapiro-Wilkův test normality, který neuvádíme). Pokud test zamítne normalitu, použijeme Wilcoxonův test, jinak použijeme dvouvýběrový t -test. Od používání takového postupu zásadně odrazujeme. Jak víme, jedná se o dva testy, které testují rozdílné hypotézy, nemůžeme je tedy použít na ten samý problém. Pokud si nejsme jisti normalitou dat, provedeme raději rovnou Welchův test, který normalitu nevyžaduje a testuje právě tu hypotézu, která byla zadána.

Poznámka. V případě shod je zapotřebí tvrzení 6.6 mírně modifikovat. Jestliže se statistika $W_{n,m}$ počítá pomocí průměrných pořadí, tak vztah (i) platí, pokud definujeme statistiku $W_{n,m}^*$ jako

$$W_{n,m}^* = \sum_{i=1}^n \sum_{j=1}^m [\mathbb{1}\{X_i < Y_j\} + \frac{1}{2}\mathbb{1}\{X_i = Y_j\}].$$

Část (ii) je pak zapotřebí opravit na

$$\frac{W_{n,m}^*}{mn} \xrightarrow{P} P[X_i < Y_j] + \frac{1}{2}P[X_i = Y_j].$$

6.5. DVOUVÝBĚROVÝ F -TEST SHODY ROZPTYLŮ

Dvouvýběrový F -test shody rozptylů* je přesný test porovnávající rozptyly dvou nezávislých náhodných výběrů za předpokladu **normálního** rozdělení.

Model: $\mathcal{F} = \{X_i \sim N(\mu_X, \sigma_X^2), Y_j \sim N(\mu_Y, \sigma_Y^2), \mu_X, \mu_Y \in \mathbb{R}, \sigma_X^2 > 0, \sigma_Y^2 > 0\}$

Testované parametry: Rozptyly $\sigma_X^2 = \text{var } X_i$ a $\sigma_Y^2 = \text{var } Y_j$.

Hypotéza a alternativa:

$$H_0 : \sigma_X^2 = \sigma_Y^2, \quad H_1 : \sigma_X^2 \neq \sigma_Y^2.$$

Testová statistika:

$$F = \frac{S_X^2}{S_Y^2},$$

kde S_X^2 je výběrový rozptyl výběru X_1, \dots, X_n a S_Y^2 je výběrový rozptyl výběru Y_1, \dots, Y_m .

Poznámka.

- Z věty 2.11 plyne, že testová statistika má za platnosti modelu a hypotézy přesně rozdělení $F_{n-1, m-1}$.

* Angl. *two-sample F-test of equality of variances*

- Hypotézu zamítneme, pokud se výběrové rozptyly příliš liší, tj. pokud je testová statistika buď moc velká nebo moc malá.

Kritický obor:

$$H_0 \text{ zamítneme} \Leftrightarrow F \leq F_{n-1, m-1}(\alpha/2) \text{ nebo } F \geq F_{n-1, m-1}(1 - \alpha/2),$$

kde $F_{n-1, m-1}(\alpha/2)$ a $F_{n-1, m-1}(1 - \alpha/2)$ jsou po řadě $(\alpha/2)$ -tý a $(1 - \alpha/2)$ -tý kvantil F -rozdělení s $n - 1$ a $m - 1$ stupni volnosti.

P-hodnota: $p = 2 \min \{1 - G(s), G(s)\}$, kde s je pozorovaná hodnota testové statistiky F a G je distribuční funkce rozdělení $F_{n-1, m-1}$.

Interval spolehlivosti pro σ_X^2/σ_Y^2 : Dle věty 2.11 platí

$$P \left[F_{n-1, m-1}(\alpha/2) < \frac{S_X^2/\sigma_X^2}{S_Y^2/\sigma_Y^2} < F_{n-1, m-1}(1 - \alpha/2) \right] = 1 - \alpha.$$

Tudíž interval spolehlivosti pro σ_X^2/σ_Y^2 je dán předpisem

$$\left(\frac{S_X^2}{S_Y^2} \frac{1}{F_{n-1, m-1}(1 - \frac{\alpha}{2})}, \frac{S_X^2}{S_Y^2} \frac{1}{F_{n-1, m-1}(\frac{\alpha}{2})} \right).$$

Poznámka. Tento test lze převést na jednostranný test: Hypotéza $H'_0 : \sigma_X^2 \leq \sigma_Y^2$ se zamítá pouze pro příliš velké hodnoty testové statistiky, kritická hodnota je $F_{m-1, n-1}(1 - \alpha)$. Hypotéza $H''_0 : \sigma_X^2 \geq \sigma_Y^2$, se zamítá pouze pro příliš malé hodnoty testové statistiky, kritická hodnota je $F_{m-1, n-1}(\alpha)$.

PORUŠENÍ PŘEDPOKLADŮ

Při **porušení předpokladu normality** tento test **nedodrhuje hladinu** ani asymptoticky. Pro sestavení testu v tomto případě by bylo zapotřebí odvodit asymptotické rozdělení testové statistiky F za hypotézy a pracovat s tímto rozdělením. Alternativně lze využít také Leveneův test*. Ten se dá použít i na porovnání více nezávislých výběrů. Je však třeba upozornit, že tento test obecně netestuje shodu rozptylů, ale trochu jiného parametru variability.

* Angl. *Levene's test*

Přípravné příklady ke zkoušce.

1. Uvažujte $X_i \sim \text{Exp}(\lambda_1)$ a $Y_j \sim \text{Exp}(\lambda_2)$. Ukažte, že X_i a Y_j splňují zobecněný model posunutí (6.7).
Nápověda. Uvažujte $g(x) = \log x$.
2. Upravte dvouvýběrový F -test o rozptylu tak, aby testoval hypotézu $H_0 : \sigma_X^2 \leq \sigma_Y^2$ proti alternativě $H_1 : \sigma_X^2 > \sigma_Y^2$. Úpravu zdůvodněte. Jak by se počítala p -hodnota takto upraveného testu?
3. Rozhodujeme se, zda zaměstnance poslat spíše na jazykový kurz firmy Analfabet nebo jazykový kurz firmy Buran. Za tímto účelem jsme náhodně vybrali 20 zaměstnanců, které jsme náhodně rozdělili mezi tyto dva kurzy (na každý kurz šlo 10 zaměstnanců). Následující tabulka zachycuje počty bodů z testu z anglického jazyka všech zaměstnanců po absolvování kurzu

kurz Analfabet	37	41	36	48	42	36	42	44	40	34
kurz Buran	38	43	43	47	52	44	41	42	42	39

Jak byste otestovali, že mezi kurzy není rozdíl.

Nejde o samotné provedení testu, ale o jeho podrobné vysvětlení (tj. definujte vhodný model pro data, nulovou a alternativní hypotézu, testovou statistiku a kritický obor).

4. Máme k dispozici údaje o mzdách 100 zaměstnanců v jedné velké pojišťovně. Dále u těchto zaměstnanců máme informaci o tom, zda vystudovali MFF UK či jinou školu. Navrhněte test (tj. definujte vhodný model pro data, nulovou a alternativní hypotézu, testovou statistiku a vzorec pro výpočet p -hodnoty), kterým bychom chtěli prokázat, že absolventi MFF UK mají spíše vyšší mzdy než absolventi jiných škol.
5. Máme dva generátory nezávislých čísel z nějakého konkrétního rozdělení. Z každého generátoru jsme získali 500 náhodných čísel. Navrhněte test (tj. definujte vhodný model pro data, nulovou a alternativní hypotézu, testovou statistiku a kritický obor), pomocí kterého bychom otestovali, že oba generátory generují náhodná čísla ze stejného rozdělení.

7. JEDNOVÝBĚROVÉ A DVOUVÝBĚROVÉ PROBLÉMY PRO BINÁRNÍ DATA

V této kapitole se budeme zabývat *binárními veličinami*, které nabývají pouze dvou hodnot.

7.1. JEDNOVÝBĚROVÝ PROBLÉM

Alternativní rozdělení je nejjednodušším modelem pro kategoriální veličinu, která nabývá pouze dvou hodnot zakódovaných jako 0 a 1. Nechť $p_X \in (0, 1)$ je pravděpodobnost, že daný jedinec je klasifikován do kategorie 1.

Nechť Y_1, \dots, Y_n je náhodný výběr z alternativního rozdělení $\text{Alt}(p_X)$ zaznamenávající klasifikaci n jedinců do kategorií 0 a 1. Označme počet jedinců klasifikovaných do skupiny 1 jako $X_n = \sum_{i=1}^n Y_i$. Tato veličina má rozdělení $\text{Bi}(n, p_X)$ (viz věta 2.3(iv)). Počet jedinců klasifikovaných do skupiny 0 je $n - X_n \sim \text{Bi}(n, 1 - p_X)$.

Nestranným a konsistentním odhadem parametru p_X je relativní četnost

$$\hat{p}_n = \frac{X_n}{n} = \frac{\sum_{i=1}^n Y_i}{n} = \bar{Y}_n.$$

Jeho vlastnosti vycházejí z vlastností průměru a jsou shrnuty ve větě 2.3.

7.1.1. CLOPPEROVA-PEARSONOVA METODA

Nejprve se budeme zabývat metodami pro sestavení intervalu spolehlivosti pro pravděpodobnost p_X a pro testování hypotéz o p_X založenými na přesném rozdělení statistiky X_n , tj. $\text{Bi}(n, p_X)$.

Uvažujme hypotézu $H_0 : p_X = p_0$ proti alternativě $H_1 : p_X \neq p_0$. Stanovme kritický obor

$$H_0 \text{ zamítneme} \Leftrightarrow X_n \leq c_L(\alpha) \text{ nebo } X_n \geq c_U(\alpha),$$

kde $c_L(\alpha)$ je největší celé číslo, které splňuje

$$P(\text{Bi}(n, p_0) \leq c_L(\alpha)) = \sum_{j=0}^{c_L(\alpha)} \binom{n}{j} p_0^j (1 - p_0)^{n-j} \leq \frac{\alpha}{2}$$

a $c_U(\alpha)$ je nejmenší celé číslo, které splňuje

$$P(\text{Bi}(n, p_0) \geq c_U(\alpha)) = \sum_{j=c_U(\alpha)}^n \binom{n}{j} p_0^j (1 - p_0)^{n-j} \leq \frac{\alpha}{2}.$$

Tento test (zvaný *Clopperův-Pearsonův*) má hladinu, jež nepřesahuje α (vzhledem k diskrétnímu rozdělení testové statistiky nelze vždy dosáhnout stanovené hladiny α). P-hodnota tohoto testu je dána vzorcem:

$$p(x_n) = 2 \min \left\{ P(\text{Bi}(n, p_0) \leq x_n), P(\text{Bi}(n, p_0) \geq x_n) \right\} = 2 \min \left\{ G_0(x_n), 1 - G_0(x_n - 1) \right\},$$

kde G_0 je distribuční funkce $\text{Bi}(n, p_0)$ a x_n je napozorovaná hodnota X_n .

Nyní řešíme úlohu *sestavení intervalu spolehlivosti* pro p_X s pravděpodobností pokrytí (alespoň) $1 - \alpha$. Podle (dualita intervalů spolehlivosti a testování), můžeme sestavit požadovaný interval spolehlivosti jako množinu obsahující všechny parametry $p \in (0, 1)$, pro něž při pozorovaných datech X_n Clopperův-Pearsonův test nezamítá hypotézu $H_0 : p_X = p$ proti alternativě $H_1 : p_X \neq p$. Označme ještě G_p distribuční funkci rozdělení $\text{Bi}(n, p)$. Potom

$$\begin{aligned} IS_n &= \{p \in (0, 1) : p(X_n) > \alpha, \text{ kde } p(X_n) \text{ je } p\text{-hodnota testu } H_0 : p_X = p\} \\ &= \{p \in (0, 1) : 2 \min\{G_p(X_n), 1 - G_p(X_n - 1)\} > \alpha\} \\ &= \left\{ p \in (0, 1) : \sum_{j=0}^{X_n} \binom{n}{j} p^j (1-p)^{n-j} > \frac{\alpha}{2} \quad \text{a zároveň} \quad \sum_{j=X_n}^n \binom{n}{j} p^j (1-p)^{n-j} > \frac{\alpha}{2} \right\}. \end{aligned}$$

Tj. interval spolehlivosti bude tvaru (p_L, p_U) , kde p_L a p_U nalezneme jako řešení následujících rovnic

$$\sum_{j=X_n}^n \binom{n}{j} p^j (1-p)^{n-j} = \frac{\alpha}{2}, \quad \sum_{j=0}^{X_n} \binom{n}{j} p^j (1-p)^{n-j} = \frac{\alpha}{2}.$$

Dá se ukázat, že p_L a p_U lze explicitně vyjádřit a dostáváme interval ve tvaru

$$\left(\frac{X_n q_L(\alpha)}{X_n q_L(\alpha) + n - X_n + 1}, \frac{(X_n + 1) q_U(\alpha)}{(X_n + 1) q_U(\alpha) + n - X_n} \right),$$

kde $q_L(\alpha)$ je $\alpha/2$ -kvantil rozdělení $F_{2X_n, 2(n-X_n+1)}$ a $q_U(\alpha)$ je $(1-\alpha/2)$ -kvantil $F_{2(X_n+1), 2(n-X_n)}$. Pokud $X_n = 0$, položíme dolní mez intervalu rovnou 0, pokud $X_n = n$, položíme horní mez intervalu rovnou 1.

Výše uvedený interval se nazývá *Clopperův-Pearsonův interval spolehlivosti* pro parametr binomického rozdělení. Výhodou tohoto intervalu je, že pravděpodobnost pokrytí má hodnotu alespoň $1 - \alpha$ pro jakýkoliv rozsah výběru. Jeho nevýhodou je, že jeho pravděpodobnost pokrytí může být o hodně vyšší než $1 - \alpha$ a že mívá příliš velkou délku.

Nyní se můžeme vrátit ke Clopperově-Pearsonově testu hypotézy $H_0 : p_X = p_0$ proti alternativě $H_1 : p_X \neq p_0$. Místo toho, abychom složitě počítali kritické hodnoty $c_L(\alpha)$ a $c_U(\alpha)$, spočítáme Clopperův-Pearsonův interval spolehlivosti a H_0 zamítneme, pokud p_0 v tomto intervalu neleží.

7.1.2. KLASICKÁ ASYMPTOTICKÁ METODA

V příkladě uvedeném v kapitole 3.4.2 na str. 56 jsme odvodili asymptotický interval spolehlivosti pro p_X založený na bodě (iii) věty 2.3 a Sluckého větě. Podle (3.6) platí

$$Z_n = \frac{\sqrt{n}(\hat{p}_n - p_X)}{\sqrt{\hat{p}_n(1 - \hat{p}_n)}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Toto tvrzení lze použít k odvození asymptotického testu hypotézy $H_0 : p_X = p_0$ proti alternativě $H_1 : p_X \neq p_0$ s kritickým oborem

$$H_0 \text{ zamítneme} \Leftrightarrow \left| \frac{\sqrt{n}(\hat{p}_n - p_0)}{\sqrt{\hat{p}_n(1 - \hat{p}_n)}} \right| \geq u_{1-\alpha/2}. \quad (7.1)$$

Z duality testování a intervalů spolehlivosti (tvrzení 4.2(ii)), můžeme interval spolehlivosti pro p_X najít jako

$$IS_n = \left\{ p \in (0, 1) : \left| \frac{\sqrt{n}(\hat{p}_n - p)}{\sqrt{\hat{p}_n(1 - \hat{p}_n)}} \right| < u_{1-\alpha/2} \right\} = \left(\hat{p}_n - u_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_n(1 - \hat{p}_n)}{n}}, \hat{p}_n + u_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_n(1 - \hat{p}_n)}{n}} \right).$$

Za povšimnutí stojí to, že tento interval spolehlivosti odpovídá intervalu spolehlivosti pro p_X z kapitoly 3.4.2.

Nevýhodou tohoto přístupu je, že pokud p_X je blízké nule nebo jedné, pak je zapotřebí relativně velký počet pozorování, aby asymptotická aproximace fungovala spolehlivě. V praxi se někdy doporučuje, že je potřeba mít alespoň 5 úspěchů a alespoň 5 neúspěchů. Za povšimnutí také stojí, že krajní body intervalu spolehlivosti mohou být menší než 0 nebo větší než 1.

Cvičení. Jelikož alternativní rozdělení náleží do \mathcal{L}_+^2 , tak bychom také mohli použít (asymptotický) t -test, viz kapitola 5.3. Ukažte, že v tomto případě by měla testová statistika tvar

$$T_n = \frac{\sqrt{n-1}(\hat{p}_n - p_0)}{\sqrt{\hat{p}_n(1 - \hat{p}_n)}}$$

a porovnávali bychom ji s kvantily t_{n-1} -rozdělení. Dostali bychom tedy test, který je o trochu konzervativnější než test uvedený v (7.1).

7.1.3. WILSONOVA METODA

Wilsonova metoda je založena přímo na bodě (iii) věty 2.3

$$W_n = \frac{\sqrt{n}(\hat{p}_n - p_X)}{\sqrt{p_X(1 - p_X)}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1)$$

bez aplikace Sluckého věty. Za platnosti hypotézy $H_0 : p_X = p_0$ známe p_X a toho využijeme k sestavení kritického oboru

$$H_0 \text{ zamítneme} \Leftrightarrow \left| \frac{\sqrt{n}(\hat{p}_n - p_0)}{\sqrt{p_0(1 - p_0)}} \right| \geq u_{1-\alpha/2}.$$

Tento test se nazývá *Wilsonův*.

Interval spolehlivosti najdeme opět pomocí duality mezi testováním a intervaly spolehlivosti jako

$$IS_n = \left\{ p \in (0, 1) : \left| \frac{\sqrt{n}(\hat{p}_n - p)}{\sqrt{p(1-p)}} \right| < u_{1-\alpha/2} \right\}.$$

Po úpravě (vyřešení kvadratické rovnici pro p) dostaneme asymptotický interval s krajními body

$$\left(\hat{p}_n + \frac{u^2}{2n} \mp u \sqrt{\frac{\hat{p}_n(1-\hat{p}_n)}{n} + \frac{u^2}{4n^2}} \right) \frac{1}{1 + u^2/n},$$

kde u je zkrácené značení pro $u_{1-\alpha/2}$. Tento interval se též nazývá *Wilsonův*. V literatuře se uvádí, že Wilsonův test a interval dává přesnější výsledky než metody z kapitoly 7.1.2.

Je zajímavé si povšimnout, že střed Wilsonova intervalu lze vyjádřit jako vážený průměr $w_n \hat{p}_n + (1 - w_n)1/2$, kde $w_n = (1 + u^2/n)^{-1} \rightarrow 1$ pro $n \rightarrow \infty$. Počítáme-li 95% interval spolehlivosti, pak střed Wilsonova intervalu je zhruba $(X_n + 2)/(n + 4)$.

7.1.4. LOGITOVÁ METODA

Logitová metoda je založena na šanci místo na pravděpodobnosti.

Definice 7.1 Nechť úspěch nastává s pravděpodobností p . Podíl $\frac{p}{1-p}$ pravděpodobnosti úspěchu a neúspěchu se nazývá *šance** na úspěch.

Pojem šance se běžně používá při kursových sázkách.

Zvolme jako odhadovaný parametr logaritmus šance $\theta_X = \log \frac{p_X}{1-p_X}$. Tomuto parametru se běžně říká *logit*, transformace $g(x) = \log \left(\frac{x}{1-x} \right)$ se nazývá *logitová*. Logitová transformace $g(x)$ je rostoucí a spojitě diferencovatelná pro $x \in (0, 1)$ a zobrazuje interval $(0, 1)$ na \mathbb{R} . Inversní transformace je $g^{-1}(y) = \frac{\exp\{y\}}{1+\exp\{y\}}$. Logaritmus šance θ_X tedy může nabývat libovolné hodnoty v \mathbb{R} a můžeme z ní vyjádřit pravděpodobnost p_X jako $p_X = \exp\{\theta_X\}/(1 + \exp\{\theta_X\})$.

Logaritmus šance θ_X odhadneme transformací $g(\hat{p}_n)$ odhadu \hat{p}_n . Dostaneme odhad

$$\hat{\theta}_n = \log \left(\frac{\hat{p}_n}{1-\hat{p}_n} \right),$$

který je podle zákona velkých čísel (tvrzení 1.5) a věty o spojitě transformaci (tvrzení 1.2) konsistentním (ne však nestranným) odhadem θ_X .

Asymptotické rozdělení $\hat{\theta}_n$ získáme aplikací bodu (iii) věty 2.3 a delta metody (tvrzení 1.7).

Věta 7.1 Nechť $p_X \in (0, 1)$. Pak platí

(i)

$$\sqrt{n}(\hat{\theta}_n - \theta_X) \xrightarrow[n \rightarrow \infty]{d} N\left(0, \frac{1}{p_X} + \frac{1}{1-p_X}\right),$$

* Angl. *odds*

(ii)

$$\sqrt{\frac{X_n(n-X_n)}{n}} (\widehat{\theta}_n - \theta_X) \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Označme $D_n = \sqrt{\frac{n}{X_n(n-X_n)}}$. Je to vlastně odhad směrodatné chyby $\widehat{\theta}_n$.

Na základě věty 7.1 můžeme sestavit asymptotický test hypotézy $H_0 : p_X = p_0$. Označme $\theta_0 = \log \frac{p_0}{1-p_0}$. Hypotézu H_0 můžeme přepsat jako $H_0 : \theta_X = \theta_0$ a zamítáme ji ve prospěch alternativy $H_1 : \theta_X \neq \theta_0$ pokud

$$\frac{1}{D_n} \left| \widehat{\theta}_n - \theta_0 \right| \geq u_{1-\alpha/2}.$$

Tento test nazveme *logitový*.

Interval spolehlivosti pro θ_X s pravděpodobností pokrytí konvergující k $1 - \alpha$ má tvar

$$\left(\widehat{\theta}_n - u_{1-\frac{\alpha}{2}} D_n, \widehat{\theta}_n + u_{1-\frac{\alpha}{2}} D_n \right).$$

Aplikujeme-li rostoucí funkci g^{-1} na oba krajní body tohoto intervalu, dostaneme asymptotický $100(1 - \alpha)$ -procentní interval spolehlivosti pro p_X ve tvaru

$$\left(\frac{\frac{\widehat{p}_n}{1-\widehat{p}_n} e^{-u_{1-\alpha/2} D_n}}{1 + \frac{\widehat{p}_n}{1-\widehat{p}_n} e^{-u_{1-\alpha/2} D_n}}, \frac{\frac{\widehat{p}_n}{1-\widehat{p}_n} e^{u_{1-\alpha/2} D_n}}{1 + \frac{\widehat{p}_n}{1-\widehat{p}_n} e^{u_{1-\alpha/2} D_n}} \right). \quad (7.2)$$

Interval (7.2) nazýváme *logitový*. Oba jeho krajní body jistě leží uvnitř $(0, 1)$. Navíc konvergence $\widehat{\theta}_n$ k normálnímu rozdělení je rychlejší než konvergence \widehat{p}_n , takže limitní aproximace založená na $\widehat{\theta}_n$ je přesnější než aproximace založená na \widehat{p}_n . Logitová metoda patří spolu s Wilsonovou k metodám doporučovaným v literatuře.

7.2. DVOUVÝBĚROVÝ PROBLÉM

Nechť Y_{11}, \dots, Y_{1n} je náhodný výběr z alternativního rozdělení $\text{Alt}(p_1)$ a Y_{21}, \dots, Y_{2m} je náhodný výběr z $\text{Alt}(p_2)$. Označme $X_1 = \sum_{i=1}^n Y_{1i}$ a $X_2 = \sum_{i=1}^m Y_{2i}$. Budeme se tedy zabývat porovnáním dvou nezávislých binomických veličin $X_1 \sim \text{Bi}(n, p_1)$ a $X_2 \sim \text{Bi}(m, p_2)$. Chceme zjistit, zdali a jakým způsobem se liší pravděpodobnosti p_1 a p_2 . Jejich odlišnost můžeme vyjádřit různými způsoby, z toho nám vyplyne několik variant odhadů a testů.

Pokud veličiny X_1 a X_2 udávají počty nějakých negativních událostí (smrt, nemoc, ztráta zaměstnání, porucha, bankrot) parametry p_1 a p_2 nazýváme *riziky* události v obou populacích. Pravděpodobnosti (rizika) p_1 a p_2 můžeme odhadnout relativními četnostmi $\widehat{p}_1 = X_1/n$, $\widehat{p}_2 = X_2/m$. Jejich vlastnosti shrnuje věta 2.3.

Pravděpodobnosti (rizika) \widehat{p}_1 a \widehat{p}_2 zpravidla porovnávané jedním ze tří následujícím způsobů:

1. rozdíl pravděpodobností (nárůst rizika)* $d_X = p_1 - p_2$, odhadujeme pomocí $\widehat{d} = \widehat{p}_1 - \widehat{p}_2$;
2. podíl pravděpodobností (relativní riziko†) $r_X = \frac{p_1}{p_2}$, odhadujeme pomocí $\widehat{r} = \frac{\widehat{p}_1}{\widehat{p}_2}$;
3. poměr šancí‡ $o_X = \frac{p_1/(1-p_1)}{p_2/(1-p_2)} = \frac{p_1(1-p_2)}{p_2(1-p_1)}$, odhadujeme pomocí $\widehat{o} = \frac{\widehat{p}_1(1-\widehat{p}_2)}{\widehat{p}_2(1-\widehat{p}_1)} = \frac{X_1(m-X_2)}{X_2(n-X_1)}$.

Pro každé z těchto porovnání budeme potřebovat asymptotické rozdělení příslušného odhadu. U všech asymptotických výsledků budeme stejně jako v kapitole 6 předpokládat, že

$$n \rightarrow \infty, \quad m \rightarrow \infty, \quad n/m \rightarrow q \in (0, \infty). \quad (7.3)$$

Výsledky uváděné v této kapitole však platí i tehdy, je-li pevný pouze celkový počet pozorování $n+m$, zatímco rozsahy výběrů n a m jsou náhodné (viz diskuse na str. 114).

Všimněme si, že pomocí centrální limitní věty máme

$$\sqrt{n}(\widehat{p}_1 - p_1) \xrightarrow{d} N(0, p_1(1-p_1)) \quad \text{a} \quad \sqrt{m}(\widehat{p}_2 - p_2) \xrightarrow{d} N(0, p_2(1-p_2)).$$

Dále díky nezávislosti \widehat{p}_1 a \widehat{p}_2 dostaneme stejným způsobem jako v důkazu věty 6.2

$$\sqrt{m} \begin{pmatrix} \widehat{p}_1 - p_1 \\ \widehat{p}_2 - p_2 \end{pmatrix} \xrightarrow{d} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \frac{p_1(1-p_1)}{q} & 0 \\ 0 & p_2(1-p_2) \end{pmatrix} \right). \quad (7.4)$$

7.2.1. ROZDÍLY PRAVDĚPODOBNOSTÍ, NÁRŮST RIZIKA

Odlišnost obou rozdělení můžeme vyjádřit např. *rozdílem pravděpodobností (rizik)* $d_X = p_1 - p_2$, jež říká, o kolik je větší riziko v populaci 1 než v populaci 2. Tento parametr může nabývat hodnot -1 až 1 , nulová hodnota odpovídá totožným pravděpodobnostem v obou populacích.

Nestranným a konsistentním odhadem parametru d_X je $\widehat{d} = \widehat{p}_1 - \widehat{p}_2$.

Tvrzení 7.2 Nechť $p_1, p_2 \in (0, 1)$ a platí (7.3). Potom

$$\frac{\widehat{d} - d_X}{\sqrt{\frac{\widehat{p}_1(1-\widehat{p}_1)}{n} + \frac{\widehat{p}_2(1-\widehat{p}_2)}{m}}} \xrightarrow{d} N(0, 1).$$

Důkaz. Budeme postupovat podobně jako v důkazu věty 6.2.

Nejdříve přepíšeme

$$\frac{\widehat{d} - d_X}{\sqrt{\frac{\widehat{p}_1(1-\widehat{p}_1)}{n} + \frac{\widehat{p}_2(1-\widehat{p}_2)}{m}}} = \frac{\sqrt{m}(\widehat{d} - d_X)}{\sqrt{\widehat{p}_1(1-\widehat{p}_1)\frac{m}{n} + \widehat{p}_2(1-\widehat{p}_2)}}.$$

* Angl. *risk difference, excess risk* † Angl. *relative risk* ‡ Angl. *odds ratio*

Nyní s pomocí zákona velkých čísel (tvrzení 1.5) a věty o spojitě transformaci (tvrzení 1.2) se ukáže, že

$$\sqrt{\widehat{p}_1(1-\widehat{p}_1)\frac{m}{n} + \widehat{p}_2(1-\widehat{p}_2)} \xrightarrow{P} \sqrt{\frac{p_1(1-p_1)}{q} + p_2(1-p_2)}.$$

Pomocí Cramérový-Sluckého věty (věta 1.3) tedy zbývá dokázat, že

$$\sqrt{m}(\widehat{d} - d_X) \xrightarrow{d} N\left(0, \frac{p_1(1-p_1)}{q} + p_2(1-p_2)\right),$$

což plyne podobně jako v důkazu věty 6.2 ze sdružené asymptotické normality odhadů \widehat{p}_1 a \widehat{p}_2 v (7.4). \square

Pro asymptotický test hypotézy $H_0 : d_X = 0$ proti alternativě $H_1 : d_X \neq 0$ použijeme testovou statistiku

$$\widetilde{T}_d = \frac{\widehat{d}}{\sqrt{\frac{\widehat{p}_1(1-\widehat{p}_1)}{n} + \frac{\widehat{p}_2(1-\widehat{p}_2)}{m}}}$$

a hypotézu zamítneme pokud $|\widetilde{T}_d| \geq u_{1-\alpha/2}$.

Z tvrzení 7.2 dostaneme postupnými úpravami

$$P\left[\widehat{d} - u_{1-\alpha/2}\sqrt{\frac{\widehat{p}_1(1-\widehat{p}_1)}{n} + \frac{\widehat{p}_2(1-\widehat{p}_2)}{m}} < d_X < \widehat{d} + u_{1-\alpha/2}\sqrt{\frac{\widehat{p}_1(1-\widehat{p}_1)}{n} + \frac{\widehat{p}_2(1-\widehat{p}_2)}{m}}\right] \rightarrow 1 - \alpha.$$

Odtud získáme asymptotický interval spolehlivosti pro rozdíl pravděpodobností d_X .

Poznámka. Jelikož za nulové hypotézy $H_0 : d_X = 0$ je $p_1 = p_2$, tak lze místo T_d použít testovou statistiku

$$T_d = \frac{\widehat{d}}{\sqrt{\widetilde{p}(1-\widetilde{p})\left(\frac{1}{n} + \frac{1}{m}\right)}}, \quad (7.5)$$

kde $\widetilde{p} = \frac{X_1+X_2}{n+m}$ je odhad společné pravděpodobnosti úspěchu za nulové hypotézy. Testová statistika T_d má za nulové hypotézy asymptoticky rozdělení $N(0, 1)$. Výhodou této statistiky je, že se ukazuje, že skutečná hladina testu založeného na T_d je zpravidla bližší předepsané hladině, než skutečná hladina testu \widetilde{T}_d . Nevýhodou však je, že statistika T_d se **nedá** využít ke konstrukci intervalu spolehlivosti pro rozdíl pravděpodobností $d_X = p_1 - p_2$.

Cvičení. Alternativně bychom mohli použít (asymptotický) dvouvýběrový t-test (viz kapitola 6.2) hypotézy $H_0 : \mu_X = \mu_Y$. Dokažte, že v tomto případě má testová statistika $\widetilde{T}_{n,m}$ tvar

$$\widetilde{T}_{n,m} = \frac{\widehat{d}}{\sqrt{\frac{\widehat{p}_1(1-\widehat{p}_1)}{n-1} + \frac{\widehat{p}_2(1-\widehat{p}_2)}{m-1}}}.$$

7.2.2. PODÍLY PRAVDĚPODOBNOSTÍ, RELATIVNÍ RIZIKO

Jiný způsob, jak vyjádřit odlišnost pravděpodobností (rizik), je *relativní riziko* $r_X = p_1/p_2$. Tento parametr říká, kolikrát je větší riziko v populaci 1 než v populaci 2 a může nabývat hodnot v intervalu $(0, \infty)$. Pravděpodobnosti (rizika) v obou populacích jsou totožné právě když $r_X = 1$.

Konsistentním (nikoli nestranným) odhadem parametru r_X je $\hat{r} = \hat{p}_1/\hat{p}_2$.

I když bychom mohli odvodit asymptotické rozdělení odhadu $\hat{r} = \hat{p}_1/\hat{p}_2$, tak se ukazuje, že normální aproximace funguje rychleji pro logaritmus tohoto podílu.

Tvrzení 7.3 Necht' $p_1, p_2 \in (0, 1)$ a platí (7.3). Potom

$$\frac{\log \hat{r} - \log r_X}{\sqrt{\frac{1-\hat{p}_1}{n\hat{p}_1} + \frac{1-\hat{p}_2}{m\hat{p}_2}}} \xrightarrow{d} N(0, 1).$$

Důkaz. Opět budeme postupovat podobně jako v důkazu věty 6.2.

Nejdříve přepíšeme

$$\frac{\log \hat{r} - \log r_X}{\sqrt{\frac{1-\hat{p}_1}{n\hat{p}_1} + \frac{1-\hat{p}_2}{m\hat{p}_2}}} = \frac{\sqrt{m} (\log \hat{r} - \log r_X)}{\sqrt{\frac{m}{n} \frac{1-\hat{p}_1}{\hat{p}_1} + \frac{1-\hat{p}_2}{\hat{p}_2}}}.$$

Nyní s pomocí zákona velkých čísel (tvrzení 1.5) a věty o spojitě transformaci (tvrzení 1.2) se ukáže, že

$$\sqrt{\frac{m}{n} \frac{1-\hat{p}_1}{\hat{p}_1} + \frac{1-\hat{p}_2}{\hat{p}_2}} \xrightarrow{P} \sqrt{\frac{1-p_1}{qp_1} + \frac{1-p_2}{p_2}}.$$

Pomocí Cramérový-Sluckého věty (věta 1.3) tedy zbývá dokázat, že

$$\sqrt{m} (\log \hat{r} - \log r_X) \xrightarrow{d} N\left(0, \frac{1-p_1}{qp_1} + \frac{1-p_2}{p_2}\right).$$

To však plyne z delta-metody (tvrzení 1.7) a ze sdružené asymptotické normality (7.4), neboť gradient funkce

$$g(p_1, p_2) = \log\left(\frac{p_1}{p_2}\right) = \log p_1 - \log p_2$$

je $(\frac{1}{p_1}, \frac{-1}{p_2})$. Tudíž asymptotický rozptyl náhodné veličiny $\sqrt{m} (\log \hat{r} - r_X)$ je

$$\begin{pmatrix} \frac{1}{p_1} & \frac{-1}{p_2} \end{pmatrix} \begin{pmatrix} \frac{p_1(1-p_1)}{q} & 0 \\ 0 & p_2(1-p_2) \end{pmatrix} \begin{pmatrix} \frac{1}{p_1} \\ \frac{-1}{p_2} \end{pmatrix} = \frac{1-p_1}{qp_1} + \frac{1-p_2}{p_2}.$$

□

Chceme otestovat, jestli $\log r_X = 0$ neboli $r_X = 1$. Pro asymptotický test hypotézy $H_0 : r_X = 1$ proti alternativě $H_1 : r_X \neq 1$ použijeme testovou statistiku

$$T_r = \frac{\log \hat{r}}{\sqrt{\frac{1-\hat{p}_1}{n\hat{p}_1} + \frac{1-\hat{p}_2}{m\hat{p}_2}}}$$

a hypotézu zamítneme pokud $|T_r| \geq u_{1-\alpha/2}$.

Z tvrzení 7.3 dostáváme

$$P \left[\log \hat{r} - u_{1-\alpha/2} \sqrt{\frac{1-\hat{p}_1}{n\hat{p}_1} + \frac{1-\hat{p}_2}{m\hat{p}_2}} < \log r_X < \log \hat{r} + u_{1-\alpha/2} \sqrt{\frac{1-\hat{p}_1}{n\hat{p}_1} + \frac{1-\hat{p}_2}{m\hat{p}_2}} \right] \rightarrow 1 - \alpha.$$

Odtud pak asymptotický interval spolehlivosti pro r_X ve tvaru

$$\left(\hat{r} \exp \left\{ -u_{1-\alpha/2} \sqrt{\frac{1-\hat{p}_1}{n\hat{p}_1} + \frac{1-\hat{p}_2}{m\hat{p}_2}} \right\}, \hat{r} \exp \left\{ u_{1-\alpha/2} \sqrt{\frac{1-\hat{p}_1}{n\hat{p}_1} + \frac{1-\hat{p}_2}{m\hat{p}_2}} \right\} \right).$$

Cvičení. Jak by vypadal kritický obor hypotézy $H_0 : r_X = 2$ proti alternativě $H_1 : r_X \neq 2$?

7.2.3. POMĚR ŠANCÍ

Třetím možným způsobem vyjádření odlišnosti dvou pravděpodobností je *poměr šancí*

$$o_X = \frac{p_1/(1-p_1)}{p_2/(1-p_2)} = \frac{p_1(1-p_2)}{p_2(1-p_1)}.$$

Tento parametr říká, kolikrát je větší šance v populaci 1 než v populaci 2. Může nabývat hodnot v intervalu $(0, \infty)$. Pravděpodobnosti (rizika) v obou populacích jsou totožná právě když $o_X = 1$.

Konsistentním (nikoli nestranným) odhadem parametru o_X je

$$\hat{o} = \frac{\hat{p}_1(1-\hat{p}_2)}{\hat{p}_2(1-\hat{p}_1)} = \frac{X_1(m-X_2)}{X_2(n-X_1)}.$$

I když bychom mohli odvodit asymptotické rozdělení odhadu $\hat{o} = \frac{\hat{p}_1(1-\hat{p}_2)}{\hat{p}_2(1-\hat{p}_1)}$, tak se ukazuje, že normální aproximace funguje rychleji pro logaritmus tohoto odhadu.

Tvrzení 7.4 Necht' $p_1, p_2 \in (0, 1)$ a platí (7.3). Položme

$$\hat{V}_o = \frac{1}{n\hat{p}_1} + \frac{1}{n(1-\hat{p}_1)} + \frac{1}{m\hat{p}_2} + \frac{1}{m(1-\hat{p}_2)} = \frac{1}{X_1} + \frac{1}{n-X_1} + \frac{1}{X_2} + \frac{1}{m-X_2}.$$

Pak

$$\frac{\log \hat{o} - \log o_X}{\sqrt{\hat{V}_o}} \xrightarrow{d} N(0, 1).$$

Důkaz. Podobně jako v důkazu věty 6.2 nejdříve přepíšeme

$$\frac{\log \hat{o} - \log o_X}{\sqrt{\hat{V}_o}} = \frac{\sqrt{m} (\log \hat{o} - \log o_X)}{\sqrt{m \hat{V}_o}}.$$

Nyní s pomocí zákona velkých čísel (tvrzení 1.5) a věty o spojité transformaci (tvrzení 1.2) se ukáže, že

$$\sqrt{m \widehat{V}_o} = \sqrt{\frac{m}{n \widehat{p}_1} + \frac{m}{n(1-\widehat{p}_1)} + \frac{1}{\widehat{p}_2} + \frac{1}{(1-\widehat{p}_2)}} \xrightarrow{P} \sqrt{\frac{1}{q p_1} + \frac{1}{q(1-p_1)} + \frac{1}{p_2} + \frac{1}{(1-p_2)}}.$$

Pomocí Cramérový-Sluckého věty (věta 1.3) tedy zbývá dokázat, že

$$\sqrt{m} (\log \widehat{o}_X - \log o_X) \xrightarrow{d} N\left(0, \frac{1}{q p_1} + \frac{1}{q(1-p_1)} + \frac{1}{p_2} + \frac{1}{(1-p_2)}\right),$$

což plyne použitím delta-metody (tvrzení 1.7) z (7.4), neboť gradient funkce

$$g(p_1, p_2) = \log\left(\frac{p_1}{1-p_1} / \frac{p_2}{1-p_2}\right) = \log p_1 - \log(1-p_1) - \log p_2 - \log(1-p_2)$$

$$\text{je } \left(\frac{1}{p_1} + \frac{1}{1-p_1}, -\frac{1}{p_2} - \frac{1}{1-p_2}\right).$$

□

Pravděpodobnosti (šance) v obou populacích jsou totožné právě když $o_X = 1$ neboli $\log o_X = 0$. Pro asymptotický test hypotézy $H_0 : o_X = 1$ proti alternativě $H_1 : o_X \neq 1$ použijeme testovou statistiku

$$T_o = \frac{\log \widehat{o}_X}{\sqrt{\widehat{V}_o}}$$

a hypotézu zamítneme pokud $|T_o| \geq u_{1-\alpha/2}$.

Z tvrzení 7.4 dostáváme

$$P\left[\log \widehat{o}_X - u_{1-\alpha/2} \sqrt{\widehat{V}_o} < \log o_X < \log \widehat{o}_X + u_{1-\alpha/2} \sqrt{\widehat{V}_o}\right] \rightarrow 1 - \alpha,$$

odkud dostáváme asymptotický interval spolehlivosti pro poměr šancí o_X ve tvaru

$$\left(\widehat{o}_X \exp\left\{-u_{1-\alpha/2} \sqrt{\widehat{V}_o}\right\}, \widehat{o}_X \exp\left\{u_{1-\alpha/2} \sqrt{\widehat{V}_o}\right\}\right).$$

Cvičení. Jak by vypadal kritický obor hypotézy $H_0 : o_X \leq 2$ proti alternativě $H_1 : o_X > 2$?

Přípravné příklady ke zkoušce.

Vaše řešení by mělo **vždy obsahovat matematický model**. V případě testování hypotéz, pak také testovou statistiku a její přesné (či asymptotické) rozdělení za nulové hypotézy. Dále pak kritický obor nebo vzorec pro výpočet p -hodnoty. Mělo by být také řečeno, zda je daný test přesný nebo asymptotický. Podobně u intervalu spolehlivosti.

1. Mezi 100 dotázanými vysokoškoláky by volilo hnutí *Autoritářsky Nacionálního Obrození* 11 respondentů, zatímco mezi 200 respondenty s nejvýše středoškolským vzděláním by to bylo 84 respondentů.
 - a) Může se být hnutí jisté tím, že mezi respondenty s nejvýše středoškolským vzděláním má podporu alespoň 35 % ?
 - b) Dá se říct, že respondent s vysokoškolským vzděláním má alespoň dvakrát menší pravděpodobnost, že bude volit dané hnutí?
2. Starosta by rád uspořádal novoroční ohňostroj, ale není si jistý, jak by to občané přijali. Proto si nechal udělat průzkum, ve kterém ze 100 dotazovaných občanů by 61 novoroční ohňostroj přivítalo. Může mít starosta na základě těchto dat dostatečnou jistotu, že více než polovina občanů by novoroční ohňostroj uvítalo?

8. MULTINOMICKÉ ROZDĚLENÍ A KONTINGENČNÍ TABULKY

V této kapitole se budeme zabývat *kategoriálními veličinami*, které mohou obecně nabývat dvou nebo více hodnot. Pojem kategoriální veličina byl vyložen v kapitole 3.2.2. Stručně řečeno, jde o diskrétní veličinu nabývající konečně mnoha hodnot, typicky $1, \dots, K$, jejíž hodnoty nemusí mít numerickou interpretaci, ale označují členství v nějaké skupině (kategorii). Parametry používané v analýze kategoriálních dat jsou typicky pravděpodobnosti jednotlivých hodnot.

8.1. MULTINOMICKÉ ROZDĚLENÍ

Multinomické rozdělení zobecňuje binomické rozdělení na situaci, kdy kategoriální veličina může nabývat více než dvou hodnot.

MULTINOMICKÉ ROZDĚLENÍ: DEFINICE A VLASTNOSTI

Definice 8.1 (Multinomické rozdělení) Nechť $K \geq 2$ a $n \geq 1$ jsou přirozená čísla a $\mathbf{p} = (p_1, \dots, p_K)^\top$ je vektor konstant splňující $p_k > 0 \forall k$ a $\sum_{k=1}^K p_k = 1$. Náhodný vektor $\mathbf{X} = (X_1, \dots, X_K)^\top$ má multinomické rozdělení $\text{Mult}_K(n, \mathbf{p})$, právě když jeho hustota vzhledem k součinové čítecí míře na \mathbb{Z}^K je

$$P[X_1 = x_1, X_2 = x_2, \dots, X_K = x_K] = \begin{cases} \frac{n!}{x_1! \cdots x_K!} p_1^{x_1} p_2^{x_2} \cdots p_K^{x_K} & \sum_{k=1}^K x_k = n \\ & x_k \in \mathbb{N}_0 \quad \forall k \\ 0 & \text{jinak.} \end{cases}$$

Multinomické rozdělení je rozdělení počtu pozorování přidělených do každé z K možných příhrádek v n nezávislých experimentech, přičemž v každém experimentu jsou pravděpodobnosti přiřazení do jednotlivých příhrádek dány složkami vektoru pravděpodobností \mathbf{p} .

Věta 8.1 (Rozklad multinomického rozdělení.) Nechť $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ jsou nezávislé náhodné vektory s rozdělením $\text{Mult}_K(1, \mathbf{p})$. Pak $\sum_{i=1}^n \mathbf{Y}_i \sim \text{Mult}_K(n, \mathbf{p})$.

Důkaz. Budeme postupovat indukcí.

Pro $n = 1$ tvrzení platí triviálně.

Předpokládejme, že věta platí pro $n - 1$. Tj. $\mathbf{X} = \sum_{i=1}^{n-1} \mathbf{Y}_i \sim \text{Mult}_K(n - 1, \mathbf{p})$. Ukážeme, že $\mathbf{X} + \mathbf{Y}_n \sim \text{Mult}_K(n, \mathbf{p})$.

Označme $\mathbf{Y}_n = (Y_{n1}, \dots, Y_{nK})^\top$ a pro $\sum_{k=1}^K x_k = n$ s využitím indukčního předpokladu počítejme

$$\begin{aligned} & \mathbb{P}[X_1 + Y_{n1} = x_1, \dots, X_K + Y_{nK} = x_K] \\ &= \sum_{k=1}^K \mathbb{P}[X_1 + Y_{n1} = x_1, \dots, X_K + Y_{nK} = x_K \mid Y_{nk} = 1] \mathbb{P}[Y_{nk} = 1] \\ &= \sum_{k=1}^K \mathbb{P}[X_k = x_k - 1, X_j = x_j, \forall j \neq k] \mathbb{P}[Y_{nk} = 1] \\ &= \sum_{k=1}^K \frac{(n-1)!}{(x_k - 1)! \prod_{j=1, j \neq k}^K x_j!} p_k^{x_k - 1} \left(\prod_{j=1, j \neq k}^K p_j^{x_j} \right) p_k \\ &= \frac{(n-1)!}{\prod_{j=1}^K x_j!} \left(\prod_{j=1}^K p_j^{x_j} \right) \sum_{k=1}^K x_k = \frac{n!}{x_1! \cdots x_K!} p_1^{x_1} \cdots p_K^{x_K}. \end{aligned}$$

□

Věta 8.2 (Vlastnosti multinomického rozdělení.) Nechť $\mathbf{X} \sim \text{Mult}_K(n, \mathbf{p})$. Pak

- (i) $X_k \sim \text{Bi}(n, p_k)$,
- (ii) $E X_k = n p_k$, $\text{var } X_k = n p_k (1 - p_k)$,
- (iii) $\text{cov}(X_j, X_k) = -n p_j p_k$, pro $j \neq k$,
- (iv) varianční matice \mathbf{X} je

$$\text{var } \mathbf{X} = n [\text{diag}(\mathbf{p}) - \mathbf{p}^{\otimes 2}],$$

kde $\text{diag}(\mathbf{p})$ je diagonální matice, která má na diagonále prvky vektoru $\mathbf{p} = (p_1, \dots, p_K)$ a $\mathbf{p}^{\otimes 2} = \mathbf{p}\mathbf{p}^\top$.

Důkaz. Dle věty 8.1 si můžeme \mathbf{X} reprezentovat jako $\mathbf{X} = \sum_{i=1}^n \mathbf{Y}_i$, kde $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ jsou nezávislé náhodné vektory s rozdělením $\text{Mult}_K(1, \mathbf{p})$.

Část (i) plyne z toho, že $X_k = \sum_{i=1}^n Y_{ik}$ a část (ii) plyne z vlastností binomického rozdělení.

Část (iii). Pomocí výše uvedené reprezentace pro $j \neq k$ počítejme

$$\begin{aligned} \text{cov}(X_j, X_k) &= \text{cov}\left(\sum_{i=1}^n Y_{ij}, \sum_{l=1}^n Y_{lk}\right) = \sum_{i=1}^n \sum_{l=1}^n \text{cov}(Y_{ij}, Y_{lk}) \\ &= \sum_{i=1}^n \text{cov}(Y_{ij}, Y_{ik}) = n \text{cov}(Y_{ij}, Y_{ik}) \\ &= n (E Y_{ij} Y_{ik} - E Y_{ij} E Y_{ik}) = -n p_j p_k, \end{aligned}$$

kde jsme využili toho, že $\text{cov}(Y_{ij}, Y_{lk}) = 0$ pro $i \neq l$ (z nezávislosti náh. vektorů \mathbf{Y}_i a \mathbf{Y}_l), $E Y_{ij} Y_{ik} = 0$ (neboť pouze jedna složka vektoru \mathbf{Y}_i je nenulová), $E Y_{ij} = p_j$ a $E Y_{ik} = p_k$.

Část (iv). Z částí (ii) a (iii) plyne

$$\text{var } \mathbf{X} = \begin{pmatrix} np_1(1-p_1) & -np_1p_2 & \dots & -np_1p_K \\ -np_2p_1 & np_2(1-p_2) & \dots & -np_2p_K \\ \dots & \dots & \dots & \dots \\ -np_Kp_1 & -np_Kp_2 & \dots & np_K(1-p_K) \end{pmatrix} = n [\text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top].$$

□

Věta 8.3 (Asymptotické vlastnosti multinomického rozdělení.)

Nechť $\mathbf{X}_n = (X_{n1}, \dots, X_{nK})^\top \sim \text{Mult}_K(n, \mathbf{p})$. Pak

(i)

$$\frac{1}{\sqrt{n}}(\mathbf{X}_n - n\mathbf{p}) \xrightarrow[n \rightarrow \infty]{d} N_K(\mathbf{0}, \text{diag}(\mathbf{p}) - \mathbf{p}^{\otimes 2}),$$

(ii)

$$\sum_{k=1}^K \frac{(X_{nk} - np_k)^2}{np_k} \xrightarrow[n \rightarrow \infty]{d} \chi_{K-1}^2.$$

Důkaz. Část (i). Díky větě 8.1 si můžeme \mathbf{X}_n reprezentovat jako $\mathbf{X}_n = \sum_{i=1}^n \mathbf{Y}_i$, kde $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ jsou nezávislé náhodné vektory s rozdělením $\text{Mult}_K(1, \mathbf{p})$. Z věty 8.2 pak víme, že

$$E \mathbf{Y}_i = \mathbf{p}, \quad \text{var } \mathbf{Y}_i = \text{diag}(\mathbf{p}) - \mathbf{p}^{\otimes 2}.$$

Tedy pomocí centrální limitní věty pro nezávislé stejně rozdělené náhodné vektory (tvrzení 1.6)

$$\frac{1}{\sqrt{n}}(\mathbf{X}_n - n\mathbf{p}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbf{Y}_i - \mathbf{p}) \xrightarrow[n \rightarrow \infty]{d} N_K(\mathbf{0}, \text{diag}(\mathbf{p}) - \mathbf{p}^{\otimes 2}).$$

Část (ii). Všimněme si, že

$$\sum_{k=1}^K \frac{(X_{nk} - np_k)^2}{np_k} = \mathbf{Z}_n^\top \mathbf{Z}_n,$$

kde

$$\mathbf{Z}_n = \frac{1}{\sqrt{n}} \text{diag}\left(\frac{1}{\sqrt{p}}\right) (\mathbf{X}_n - n\mathbf{p}).$$

S využitím části (i)

$$\mathbf{Z}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{Z} \sim N_K(\mathbf{0}, \Sigma), \quad (8.1)$$

kde $\text{diag}\left(\frac{1}{\sqrt{p}}\right)$ je diagonální matice s prvky $\frac{1}{\sqrt{p_1}}, \dots, \frac{1}{\sqrt{p_K}}$ na diagonále.

$$\Sigma = \text{diag}\left(\frac{1}{\sqrt{p}}\right) [\text{diag}(\mathbf{p}) - \mathbf{p}^{\otimes 2}] \text{diag}\left(\frac{1}{\sqrt{p}}\right) = \mathbb{1}_K - \sqrt{\mathbf{p}}^{\otimes 2}.$$

Všimněme si, že

$$\begin{aligned} (\mathbb{I}_K - \sqrt{\mathbf{p}}^{\otimes 2})(\mathbb{I}_K - \sqrt{\mathbf{p}}^{\otimes 2}) &= \mathbb{I}_K - 2\sqrt{\mathbf{p}}^{\otimes 2} + \sqrt{\mathbf{p}}\sqrt{\mathbf{p}}^T\sqrt{\mathbf{p}}\sqrt{\mathbf{p}}^T \\ &= \mathbb{I}_K - 2\sqrt{\mathbf{p}}^{\otimes 2} + \sqrt{\mathbf{p}}\sqrt{\mathbf{p}}^T = \mathbb{I}_K - \sqrt{\mathbf{p}}^{\otimes 2}, \end{aligned}$$

neboť $\sqrt{\mathbf{p}}^T\sqrt{\mathbf{p}} = 1$. Tudíž matice $\mathbb{I}_K - \sqrt{\mathbf{p}}^{\otimes 2}$ je idempotentní.

Dále z (8.1) a věty o spojitě transformaci (tvrzení 1.6) víme, že

$$\mathbf{Z}_n^T \mathbf{Z}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{Z}^T \mathbf{Z}.$$

Neboť matice $\Sigma = \mathbb{I}_K - \sqrt{\mathbf{p}}^{\otimes 2}$ je idempotentní, tak použitím lemmatu A.4 s $\mathbb{A} = \mathbb{I}_K$ dostáváme, že kvadratická forma $\mathbf{Z}^T \mathbf{Z}$ má χ^2 -rozdělení s počtem stupňů volnosti

$$\text{tr}(\mathbb{A}\Sigma) = \text{tr}(\mathbb{I}_K - \sqrt{\mathbf{p}}^{\otimes 2}) = K - \sum_{k=1}^K p_k = K - 1.$$

□

8.1.1. ODHADY PARAMETRŮ MULTINOMICKÉHO ROZDĚLENÍ

Pro odhadování jednotlivých parametrů p_k , testování hypotéz o p_k a konstrukci intervalových odhadů pro p_k můžeme použít metody popsané v kapitole 7.1, neboť podle věty 8.2(i) platí $X_k \sim \text{Bi}(n, p_k)$.

Celý vektor \mathbf{p} odhadneme pomocí $\widehat{\mathbf{p}}_n = \frac{\mathbf{X}}{n}$. Sdružené asymptotické rozdělení odhadu $\widehat{\mathbf{p}}_n$ získáme z věty 8.3(i):

$$\sqrt{n}(\widehat{\mathbf{p}}_n - \mathbf{p}) = \frac{1}{\sqrt{n}}(\mathbf{X} - n\mathbf{p}) \xrightarrow[n \rightarrow \infty]{d} N_K(\mathbf{0}, \text{diag}(\mathbf{p}) - \mathbf{p}^{\otimes 2}).$$

Pro libovolný vektor konstant \mathbf{c} o délce K , platí

$$\sqrt{n}(\mathbf{c}^T \widehat{\mathbf{p}}_n - \mathbf{c}^T \mathbf{p}) \xrightarrow[n \rightarrow \infty]{d} N(0, \mathbf{c}^T [\text{diag}(\mathbf{p}) - \mathbf{p}^{\otimes 2}] \mathbf{c}).$$

Neznámý asymptotický rozptyl $V_c = \mathbf{c}^T [\text{diag}(\mathbf{p}) - \mathbf{p}^{\otimes 2}] \mathbf{c}$ můžeme odhadnout pomocí

$$\widehat{V}_c = \mathbf{c}^T [\text{diag}(\widehat{\mathbf{p}}_n) - \widehat{\mathbf{p}}_n^{\otimes 2}] \mathbf{c}.$$

Pokud $V_c \neq 0$, pak dostaneme ze Sluckého věty (tvrzení 1.3)

$$\frac{\sqrt{n}(\mathbf{c}^T \widehat{\mathbf{p}}_n - \mathbf{c}^T \mathbf{p})}{\sqrt{\widehat{V}_c}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1). \quad (8.2)$$

Odtud můžeme snadno odvodit asymptotický test hypotéz

$$H_0 : \mathbf{c}^T \mathbf{p} = \gamma_0, \quad H_1 : \mathbf{c}^T \mathbf{p} \neq \gamma_0.$$

Vezmeme testovou statistiku

$$T_c = \frac{\sqrt{n} (\mathbf{c}^\top \hat{\mathbf{p}}_n - \gamma_0)}{\sqrt{\hat{V}_c}},$$

kteřá má podle (8.2) za platnosti hypotézy asymptoticky normované normální rozdělení a H_0 zamítneme právě když $|T_c| \geq u_{1-\alpha/2}$.

Asymptotický interval spolehlivosti pro $\mathbf{c}^\top \mathbf{p}$ založený na konvergenci (8.2) jest

$$\left(\mathbf{c}^\top \hat{\mathbf{p}}_n - u_{1-\alpha/2} \sqrt{\frac{\hat{V}_c}{n}}, \mathbf{c}^\top \hat{\mathbf{p}}_n + u_{1-\alpha/2} \sqrt{\frac{\hat{V}_c}{n}} \right).$$

Vektor \mathbf{c} vybereme tak, aby součin $\mathbf{c}^\top \mathbf{p}$ vytvořil lineární kombinaci parametrů, která nás v dané aplikaci zajímá. Chceme-li například vědět, zdali pravděpodobnosti první a poslední kategorie jsou stejné, a sestavit interval spolehlivosti pro rozdíl jejich hodnot, zvolíme $\mathbf{c} = (1, 0, \dots, 0, -1)^\top$ a $\gamma_0 = 0$.

Poznámka. Podobně pomocí Δ -metody (tvrzení 1.7) lze provádět statistickou inferenci pro jakoukoliv (dostatečně hladkou) jednorozměrnou parametrickou funkci $\theta = g(\mathbf{p})$ takovou, že $V(\mathbf{p}) = \nabla g(\mathbf{p}) [\text{diag}(\mathbf{p}) - \mathbf{p}^{\otimes 2}] [\nabla g(\mathbf{p})]^\top > 0$.

8.1.2. χ^2 -TEST DOBRÉ SHODY PRO MULTINOMICKÉ ROZDĚLENÍ

Pojmem χ^2 -test dobré shody* rozumíme test hypotézy $H_0 : \mathbf{p} = \mathbf{p}^0$ založený na větě 8.3(ii). Tato hypotéza říká, že pravděpodobnosti kategorií $\mathbf{p} = (p_1, \dots, p_K)^\top$ jsou rovny předem stanoveným hypotetickým pravděpodobnostem $\mathbf{p}^0 = (p_1^0, \dots, p_K^0)^\top$, tj. $p_k = p_k^0$ pro všechna $k \in \{1, \dots, K\}$.

Platí-li hypotéza H_0 , pak testová statistika

$$\chi^2 = \sum_{k=1}^K \frac{(X_k - np_k^0)^2}{np_k^0}$$

má podle věty 8.3(ii) asymptotické rozdělení χ_{K-1}^2 .

Testová statistika porovnává pozorovanou četnost X_k v kategorii k s četností np_k^0 očekávanou za platnosti hypotézy. **Velké hodnoty** testové statistiky svědčí proti H_0 . Hypotézu H_0 zamítneme, pokud

$$H_0 \text{ zamítneme} \Leftrightarrow \chi^2 \geq \chi_{K-1}^2(1-\alpha), \quad (8.3)$$

kde $\chi_{K-1}^2(1-\alpha)$ značí $(1-\alpha)$ -kvantil rozdělení χ_{K-1}^2 .

Nechť s_x je napozorovaná hodnota testové statistiky χ^2 . Asymptotickou p-hodnota tohoto testu spočteme pomocí vzorce (4.11)

$$p(x) = 1 - G_{K-1}(s_x),$$

kde G_{K-1} je distribuční funkce χ^2 -rozdělení o $K-1$ stupních volnosti.

* Angl. χ^2 -test of goodness of fit

Poznámka. Asymptotická aproximace χ^2 rozdělením vyžaduje, aby celkový počet pozorování n byl dostatečně velký. Jako jednoduché orientační pravidlo můžeme vzít např. požadavek, aby očekávané četnosti np_k^0 překročily 5 ve všech kategoriích $k \in \{1, \dots, K\}$ *. Vyskytují-li se v hodnotách X velmi malé četnosti nebo nuly, χ^2 aproximace může být velmi nepřesná.

Poznámka. Vezmeme-li $K = 2$, $p_1^0 \equiv p_0$, $X_2 = n - X_1$, $p_2^0 = 1 - p_0$, dostaneme

$$\chi^2 = \frac{(X_1 - np_0)^2}{np_0} + \frac{[n - X_1 - n(1 - p_0)]^2}{n(1 - p_0)} = \left[\frac{\sqrt{n}(\hat{p}_n - p_0)}{\sqrt{p_0(1 - p_0)}} \right]^2, \quad \text{kde } \hat{p}_n = \frac{X_1}{n}.$$

Tedy testová statistika χ^2 testu pro $K = 2$ kategorie je rovna čtverci Wilsonovy testové statistiky uvedené v kapitole 7.1.3.

Poznámka. Za povšimnutí stojí, že pro $K > 2$ se již hypotéza a alternativa nedají vyjádřit pomocí jednorozměrného parametru. Tudíž již v tomto případě nelze nijak jednoduše využít duality intervalových odhadů a testování hypotéz (viz tvrzení 4.2). Podobný postřeh platí pro všechny následující testy (s výjimkou kapitolky 8.2.1) v této kapitole. Proto ve zbytku kapitoly již nejsou uvedeny žádné intervaly spolehlivosti.

Příklad (Je kostka pravidelná?). Hodíme n -krát kostkou a zaznamenáme, kolikrát padly výsledky 1–6: dostaneme četnosti X_1, \dots, X_6 . Nastavíme $p_k^0 = 1/6$, $k = 1, \dots, 6$. Zamítne-li χ^2 test hypotézu H_0 , prokázali jsme, že na kostce nepadají všechna čísla stejně často.

Příklad (Rodí se děti během roku rovnoměrně?). Máme dány počty dětí narozených v jednotlivých měsících během kalendářního roku: X_1, \dots, X_{12} . Nastavíme $p_k^0 = \frac{m_k}{365}$, kde m_k je počet dní v měsíci k . Zamítne-li χ^2 test hypotézu H_0 , prokázali jsme, že děti se nerodí během roku rovnoměrně.

Příklad (Pochází náhodný výběr z distribuční funkce F_0 ?). Mějme náhodný výběr Z_1, \dots, Z_n . Zajímá nás, zdali pochází z rozdělení s distribuční funkcí $F_0(x) = F(x; \theta_0)$, kde θ_0 je známo.

Stanovíme si intervaly (a_{k-1}, a_k) , $k = 1, \dots, K$, $a_0 = -\infty$, $a_K = \infty$ tak, že jejich počet K je výrazně menší než n a do každého z intervalů padne dostatečný počet pozorování. Spočítáme, kolik pozorování padlo do k -tého intervalu: $X_k = \sum_{i=1}^n \mathbb{1}_{(a_{k-1}, a_k)}(Z_i)$. Pochází-li náhodný výběr Z_1, \dots, Z_n z rozdělení s distribuční funkcí $F_0(x) = F(x; \theta_0)$, potom vektor $\mathbf{X} = (X_1, \dots, X_K)^\top$ má multinomické rozdělení $\text{Mult}_K(n, \mathbf{p}^0)$, kde pravděpodobnosti jednotlivých kategorií jsou $p_k^0 = F(a_k; \theta_0) - F(a_{k-1}; \theta_0)$.

Provedeme test hypotézy $H_0 : \mathbf{p} = \mathbf{p}^0$ testem dobré shody podle vzorce (8.3). Zamítne-li test hypotézu H_0 , prokázali jsme, že náhodný výběr Z_1, \dots, Z_n nepochází z rozdělení $F(x; \theta_0)$.

* Výpočetní prostředí R používá ještě o něco slabší podmínku, která vychází z článku Cochran (1954)

8.1.3. χ^2 -TEST DOBRÉ SHODY PRO MULTINOMICKÉ ROZDĚLENÍ S ODHADNUTÝMI PARAMETRY

Jak jsme viděli v předchozím příkladě, pravděpodobnosti kategorií p_k^0 mohou záviset na vektoru parametrů θ_0 . Test dobré shody můžeme provést podle vzorce (8.3) jen tehdy, pokud tyto parametry známe. V praxi je ovšem někdy neznáme, můžeme je nanejvýš odhadnout. Nyní si ukážeme, jak upravit test dobré shody pro takové případy.

Uvažujme *model* \mathcal{F}_0 : Nechť náhodný vektor $\mathbf{X} = (X_1, \dots, X_K)^\top$ má multinomické rozdělení $\text{Mult}_K(n, \mathbf{p}(\theta_X))$, kde $\theta_X \in \Theta \subset \mathbb{R}^d$ je neznámý d -rozměrný parametr, $d < K - 1$, a \mathbf{p} je funkce zobrazující Θ do $(0, 1)^K$ taková, že $\mathbf{p}(\theta)^\top \mathbf{1}_K = 1$ pro všechna $\theta \in \Theta$ (součet všech složek $\mathbf{p}(\theta)$ je vždy 1). Zajímá nás, zdali rozdělení \mathbf{X} lze popsat tímto modelem nebo ne.

Příklad. V nějaké populaci se určitý gen vyskytuje ve dvou variantách (alelách) A (např. tmavé oči) a a (např. světlé oči). Mezi všemi geny v celé populaci tvoří alela A podíl $\theta_X \in (0, 1)$ a alela a má podíl $1 - \theta_X$. Každý jedinec má dva exempláře příslušného genu (jeden po otci, jeden po matce). Pokud se geny míchají nezávisle (platí tzv. Hardyho-Weinbergovo ekvilibrium), pravděpodobnosti tří možných variant genotypu jedince jsou:

Genotyp	Pravděpodobnost
AA	θ_X^2
Aa	$2\theta_X(1 - \theta_X)$
aa	$(1 - \theta_X)^2$

Pozorujeme genotypy n nezávislých jedinců a označíme X_1, X_2, X_3 počty jedinců s genotypem (po řadě) AA, Aa, aa . Platí-li Hardyho-Weinbergovo ekvilibrium, pak vektor $\mathbf{X} = (X_1, X_2, X_3)^\top$ má rozdělení $\text{Mult}_3(n, \mathbf{p}(\theta_X))$, kde $\mathbf{p}(\theta_X) = (\theta_X^2, 2\theta_X(1 - \theta_X), (1 - \theta_X)^2)^\top$. Na základě pozorování \mathbf{X} chceme otestovat, zdali se populace nachází v Hardyho-Weinbergově ekvilibriu.

Parametry θ_X potřebujeme odhadnout. K tomu je přirozené využít metodu maximální věrohodnosti. Všimněme si, že logaritmická věrohodnost má tvar

$$\ell_n(\theta) = \log \left(\frac{n!}{X_1! \dots X_K!} [p_1(\theta)]^{X_1} \dots [p_K(\theta)]^{X_K} \right) = \sum_{k=1}^K X_k \log p_k(\theta) + \log \left(\frac{n!}{X_1! \dots X_K!} \right).$$

Tedy soustava věrohodnostních rovnic $\frac{\partial \ell_n(\theta)}{\partial \theta} \Big|_{\theta=\hat{\theta}_n} = \mathbf{0}$, vede k soustavě d rovnic o d neznámých $\hat{\theta}_n$.*

$$\sum_{k=1}^K \frac{X_k}{p_k(\hat{\theta}_n)} \frac{\partial p_k(\hat{\theta}_n)}{\partial \theta} = \mathbf{0}. \quad (8.4)$$

* V knihách prof. Anděla (např. v [Anděl, 1998](#)) je tato soustava rovnic odvozena jiným způsobem a metoda odhadu se nazývá *modifikovanou metodou minimálního χ^2* .

Uvažujme testování hypotézy

$$H_0 : \exists \theta_X \in \Theta \quad \mathbf{p} = \mathbf{p}(\theta_X) \quad (\text{model } \mathcal{F}_0 \text{ platí})$$

proti alternativě

$$H_1 : \forall \theta_X \in \Theta \quad \mathbf{p} \neq \mathbf{p}(\theta_X) \quad (\text{model } \mathcal{F}_0 \text{ neplatí}).$$

Nejprve získáme odhad $\hat{\theta}_n$ parametru θ_X vyřešením soustavy (8.4). Poté můžeme otestovat hypotézu H_0 testem dobré shody s odhadnutými parametry namísto parametrů skutečných. Rozdělení testové statistiky je stále za nulové hypotézy asymptoticky χ^2 , ale za každý odhadovaný parametr se ztrácí jeden stupeň volnosti.

Tvrzení 8.4 Platí-li hypotéza H_0 , pak testová statistika

$$\chi^2 = \sum_{k=1}^K \frac{[X_k - np_k(\hat{\theta}_n)]^2}{np_k(\hat{\theta}_n)}$$

má (za jistých předpokladů regularity*) asymptoticky rozdělení χ_{K-d-1}^2 , kde d je počet odhadovaných parametrů.

Platnost tohoto tvrzení plyne z teorie maximální věrohodnosti, která bude vysvětlena v navazující přednášce.

Všimněme si, že za nulové hypotézy $E X_k = np_k(\theta_X)$. Testová statistika tedy porovnává pozorovanou četnost X_k v kategorii k s $np_k(\hat{\theta}_n)$, což je vlastně odhad očekávané četnosti za platnosti hypotézy. Jelikož proti H_0 svědčí **velké hodnoty** testové statistiky, tak

$$H_0 \text{ zamítáme} \Leftrightarrow \chi^2 \geq \chi_{K-d-1}^2(1-\alpha), \quad (8.5)$$

kde $\chi_{K-d-1}^2(1-\alpha)$ značí $(1-\alpha)$ -kvantil rozdělení χ_{K-d-1}^2 .

Poznámka. I zde je nutné mít dostatečně velký počet pozorování v každé složce vektoru \mathbf{X} .

Příklad (Pochází náhodný výběr z dané parametrické rodiny rozdělení?). Mějme náhodný výběr Z_1, \dots, Z_n . Zajímá nás, zdali pochází z rozdělení $F_X(x) = F(x; \theta_X)$, kde $\theta_X \in \Theta$ není známo (např. nějaké normální, gama nebo Poissonovo rozdělení).

Stanovíme si intervaly (a_{k-1}, a_k) , $k = 1, \dots, K$, $a_0 = -\infty$, $a_K = \infty$ tak, že jejich počet K je výrazně menší než n a do každého z intervalů padne dostatečný počet pozorování. Spočítáme, kolik pozorování padlo do k -tého intervalu: $X_k = \sum_{i=1}^n \mathbb{1}_{(a_{k-1}, a_k)}(Z_i)$.

Pochází-li náhodný výběr Z_1, \dots, Z_n z rozdělení s distribuční funkcí $F(x; \theta_X)$, potom vektor $\mathbf{X} = (X_1, \dots, X_K)^\top$ má multinomické rozdělení $\text{Mult}_K(n, \mathbf{p}(\theta_X))$, kde pravděpodobnosti jednotlivých kategorií jsou $p_k(\theta_X) = F(a_k; \theta_X) - F(a_{k-1}; \theta_X)$.

* Předpoklady regularity lze najít například ve větě 12.9 v Anděl (2002). Důkaz pak v kapitole XI.3 v Anděl (1985).

Řešením soustavy (8.4) získáme odhad $\hat{\theta}_n$ parametru θ_X . Provedeme test hypotézy H_0 testem dobré shody podle vzorce (8.5). Zamítne-li test hypotézu, prokázali jsme, že náhodný výběr Z_1, \dots, Z_n nepochází z dané rodiny rozdělení.

Zdůrazněme, že k platnosti tvrzení 8.4 je zapotřebí odhadovat parametru θ_X metodou maximální věrohodnosti v modelu $X \sim \text{Mult}_K(n, \mathbf{p}(\theta_X))$. Tvrzení 8.4 by **neplatilo**, pokud bychom použili odhad metodou maximální věrohodnosti v modelu $Z_i \sim F(\cdot; \theta_X)$, tj. pro

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log f(Z_i; \theta),$$

kde $f(\cdot; \theta)$ je hustota náhodné veličiny Z_i vzhledem k nějaké σ -konečné míře μ .

8.2. KONTINGENČNÍ TABULKY

Nechť $\begin{pmatrix} X \\ Z \end{pmatrix}$ je náhodný vektor, jehož obě dvě složky jsou kategoriální. Konkrétně předpokládejme, že $X \in \{1, \dots, J\}$ a $Z \in \{1, \dots, K\}$. Nechť

$$\begin{pmatrix} X_1 \\ Z_1 \end{pmatrix}, \dots, \begin{pmatrix} X_N \\ Z_N \end{pmatrix}$$

je náhodný výběr z rozdělení daného vektorem $\begin{pmatrix} X \\ Z \end{pmatrix}$ o (pevném) rozsahu N . Označme počet jedinců klasifikovaných do j -té kategorie veličiny X a k -té kategorie veličiny Z jako

$$n_{jk} = \sum_{i=1}^N \mathbb{1}\{X_i = j, Z_i = k\}, \quad j \in \{1, \dots, J\}, k \in \{1, \dots, K\}.$$

Náhodnou veličinu n_{jk} nazýváme *pozorovanou četností** pro kombinaci kategorií j a k . Označme $p_{jk} = P[X = j, Z = k]$ a $\mathbf{p} = (p_{11}, \dots, p_{JK})^\top$. Vzhledem k tomu, že pozorované četnosti byly vyvořeny klasifikací N nezávislých jedinců do JK kategorií, náhodný vektor $\mathbf{n} = (n_{11}, \dots, n_{JK})^\top$ musí mít multinomické rozdělení $\text{Mult}_{JK}(N, \mathbf{p})$. Protože pracujeme s multinomickým rozdělením, můžeme používat všechny výsledky z kapitoly 8.1.

Označme dále

$$\begin{aligned} n_{j+} &= \sum_{k=1}^K n_{jk}, & n_{+k} &= \sum_{j=1}^J n_{jk}, & n_{++} &= \sum_{j=1}^J \sum_{k=1}^K n_{jk} = N, \\ p_{j+} &= \sum_{k=1}^K p_{jk}, & p_{+k} &= \sum_{j=1}^J p_{jk}, & p_{++} &= \sum_{j=1}^J \sum_{k=1}^K p_{jk} = 1. \end{aligned}$$

Pravděpodobnosti p_{jk} určují sdružené rozdělení X a Z , pravděpodobnosti $p_{j+} = P[X = j]$ určují marginální rozdělení X , pravděpodobnosti $p_{+k} = P[Z = k]$ určují marginální rozdělení Z .

* Angl. *observed frequency*

8. Multinomické rozdělení a kontingenční tabulky

Pozorované četnosti můžeme sestavit do tabulky, kterou nazýváme *kontingenční tabulka**.

	$Z = 1$...	$Z = K$	Σ
$X = 1$	n_{11}	...	n_{1K}	n_{1+}
$X = 2$	n_{21}	...	n_{2K}	n_{2+}
...
$X = J$	n_{J1}	...	n_{JK}	n_{J+}
Σ	n_{+1}	...	n_{+K}	N

Podobně můžeme sestavit tabulku pravděpodobností, která popisuje sdružené rozdělení vektoru $(X, Z)^T$ i marginální rozdělení veličin X a Z .

	$Z = 1$...	$Z = K$	Σ
$X = 1$	p_{11}	...	p_{1K}	p_{1+}
$X = 2$	p_{21}	...	p_{2K}	p_{2+}
...
$X = J$	p_{J1}	...	p_{JK}	p_{J+}
Σ	p_{+1}	...	p_{+K}	1

Označme ještě podmíněné pravděpodobnosti

$$P[X = j | Z = k] = p_{j(k)} = \frac{p_{jk}}{p_{+k}},$$

$$P[Z = k | X = j] = p_{(j)k} = \frac{p_{jk}}{p_{j+}}.$$

TESTOVÁNÍ NEZÁVISLOSTI χ^2 -TESTEM

Náhodné veličiny X a Z jsou nezávislé, právě když pro každé $j \in \{1, \dots, J\}$ a $k \in \{1, \dots, K\}$ platí

$$P[X = j, Z = k] = P[X = j] P[Z = k] \quad \text{neboli} \quad p_{jk} = p_{j+} p_{+k}.$$

Pokud platí hypotéza, že X a Z jsou nezávislé náhodné veličiny, pravděpodobnosti $\mathbf{p} = (p_{11}, \dots, p_{JK})^T$ specifikující multinomické rozdělení vektoru \mathbf{n} jsou funkcemi $d = J + K - 2$ parametrů

$$\boldsymbol{\theta}_X = (p_{1+}, \dots, p_{(J-1)+}, p_{+1}, \dots, p_{+(K-1)})^T.$$

Maximálně věrohodný odhad parametru $\boldsymbol{\theta}_X$ za hypotézy nezávislosti nalezneme jako řešení soustavy rovnic (8.4), které v tomto případě mají tvar

$$\sum_{j=1}^J \sum_{k=1}^K \frac{n_{jk}}{p_{jk}(\hat{\boldsymbol{\theta}}_n)} \frac{\partial p_{jk}(\hat{\boldsymbol{\theta}}_n)}{\partial \boldsymbol{\theta}} = \mathbf{0}.$$

* Angl. *contingency table*

Všimněme si, že derivujeme-li podle parametru p_{j+} dostáváme

$$\sum_{j=1}^J \sum_{k=1}^K \frac{n_{jk}}{p_{jk}(\boldsymbol{\theta})} \frac{\partial p_{jk}(\boldsymbol{\theta})}{\partial p_{j+}} = \sum_{k=1}^K \frac{n_{jk}}{p_{j+}p_{+k}} p_{+k} - \sum_{k=1}^K \frac{n_{jk}}{p_{j+}p_{+k}} p_{+k} = \sum_{k=1}^K \left(\frac{n_{jk}}{p_{j+}} - \frac{n_{jk}}{p_{j+}} \right) = \frac{n_{j+}}{p_{j+}} - \frac{n_{j+}}{p_{j+}}.$$

Dostáváme tudíž rovnice

$$\frac{n_{j+}}{\widehat{p}_{j+}} = \frac{n_{J+}}{\widehat{p}_{J+}}, \quad j = 1, \dots, J-1.$$

A analogicky pro derivace podle parametru p_{+k}

$$\frac{n_{+k}}{\widehat{p}_{+k}} = \frac{n_{+K}}{\widehat{p}_{+K}} \quad k = 1, \dots, K-1.$$

Řešením této soustavy $(J+K-2)$ rovnic (o $J+K-2$ neznámých) pak je

$$\widehat{\boldsymbol{\theta}}_n = (\widehat{p}_{1+}, \dots, \widehat{p}_{(J-1)+}, \widehat{p}_{+1}, \dots, \widehat{p}_{+(K-1)})^\top = \left(\frac{n_{1+}}{N}, \dots, \frac{n_{(J-1)+}}{N}, \frac{n_{+1}}{N}, \dots, \frac{n_{+(K-1)}}{N} \right)^\top.$$

Maximálně věrohodné odhady složek vektoru \boldsymbol{p} za hypotézy nezávislosti pak vyjdou

$$p_{jk}(\widehat{\boldsymbol{\theta}}_n) = \widehat{p}_{j+}\widehat{p}_{+k} = \frac{n_{j+}n_{+k}}{N^2}, \quad j \in \{1, \dots, J\}, k = \{1, \dots, K\}.$$

Odhadnuté očekávané četnosti v kontingenční tabulce za platnosti hypotézy jsou tedy

$$Np_{jk}(\widehat{\boldsymbol{\theta}}_n) = N\widehat{p}_{j+}\widehat{p}_{+k} = \frac{n_{j+}n_{+k}}{N}.$$

Testová statistika z tvrzení 8.4 má tedy v tomto speciálním případě pro test nezávislosti tvar

$$\chi^2 = \sum_{j=1}^J \sum_{k=1}^K \frac{\left(n_{jk} - \frac{n_{j+}n_{+k}}{N} \right)^2}{\frac{n_{j+}n_{+k}}{N}}. \quad (8.6)$$

Podle tvrzení 8.4 má tato statistika za platnosti hypotézy nezávislosti asymptoticky χ^2 -rozdělení s počtem stupňů volnosti $JK - d - 1$, kde $d = J + K - 2$. tj. $\chi_{(J-1)(K-1)}^2$. Hypotézu nezávislosti zamítneme, pokud $\chi^2 \geq \chi_{(J-1)(K-1)}^2(1 - \alpha)$.

Poznámka. Zde popsaný test se nazývá χ^2 -test nezávislosti v kontingenční tabulce. Výše uvedené můžeme shrnout následovně.

Model (pro konting. tabulku): $(\begin{smallmatrix} X_1 \\ Z_1 \end{smallmatrix}), \dots, (\begin{smallmatrix} X_N \\ Z_N \end{smallmatrix})$ je náhodný výběr z rozdělení, které je dané diskretním náhodným vektorem $(\begin{smallmatrix} X \\ Z \end{smallmatrix})$. Tedy pro četnosti v kontingenční tabulce platí

$$\boldsymbol{n} \sim \text{Mult}_{JK}(N, (p_{11}, \dots, p_{JK})), \quad \text{kde } p_{jk} = P[X = j, Z = k]. \quad (8.7)$$

Tento model budeme nazývat *sduženě multinomický model*.

Hypotéza a alternativa:

$$H_0 : X \text{ a } Z \text{ jsou nezávislé, neboli } p_{jk} = p_{j+}p_{+k} \quad \forall j \in \{1, \dots, J\}, \quad \forall k \in \{1, \dots, K\} \quad (8.8)$$

$$H_1 : X \text{ a } Z \text{ nejsou nezávislé, neboli } \exists_{j \in \{1, \dots, J\}} \exists_{k \in \{1, \dots, K\}} p_{jk} \neq p_{j+}p_{+k}$$

Testová statistika: χ^2 daná vzorcem (8.6)

Rozdělení testové statistiky za H_0 : $\chi^2 \stackrel{\text{as.}}{\sim} \chi_{(J-1)(K-1)}^2$

Kritický obor (asympt.): H_0 zamítneme $\Leftrightarrow \chi^2 \geq \chi_{(J-1)(K-1)}^2(1 - \alpha)$.

χ^2 -TEST JAKO TEST HOMOGENITY MULTINOMICKÝCH ROZDĚLENÍ

Někdy je přirozené na kontingenční tabulku nahlížet po sloupcích jako na realizace K nezávislých multinomických rozdělení. Označme $\mathbf{n}_{(k)} = (n_{1k}, \dots, n_{Jk})^T$ vektor četností v k -tém sloupci a $\mathbf{p}_{(k)} = (p_{1(k)}, \dots, p_{J(k)})^T$ vektor odpovídajících pravděpodobností, kde na $p_{j(k)}$ můžeme nahlížet jako na $p_{j(k)} = P[X = j \mid Z = k]$.

Model (pro konting. tabulku):

$$\mathbf{n}_{(k)} \sim \text{Mult}_J(n_{+k}, \mathbf{p}_{(k)}), \quad k \in \{1, \dots, K\}, \quad \text{kde } \mathbf{n}_{(1)}, \dots, \mathbf{n}_{(K)} \text{ jsou nezávislé.} \quad (8.9)$$

Tento model budeme nazývat *sloupcově multinomický model*.

Hypotéza a alternativa:

$$H_0 : \mathbf{p}_{(1)} = \dots = \mathbf{p}_{(K)}, \quad H_1 : \exists_{k, l \in \{1, \dots, K\}} \mathbf{p}_{(k)} \neq \mathbf{p}_{(l)}, \quad (8.10)$$

Nyní si uvědomme, že složky diskrétního náhodného vektoru $\begin{pmatrix} X \\ Z \end{pmatrix}$ jsou nezávislé, právě když pro všechna $j \in \{1, \dots, J\}$ a $k \in \{1, \dots, K\}$ platí

$$P[X = j \mid Z = k] = P[X = j] \quad \text{neboli} \quad p_{j(k)} = p_{j+}.$$

Tj. nulová hypotéza nezávislosti platí, právě když

$$p_{j(1)} = p_{j(2)} = \dots = p_{j(K)} \quad \text{pro všechna } j \in \{1, \dots, J\}.$$

Z předchozího plyne, že nezávislost X a Z je ekvivalentní tomu, že vektory podmíněných pravděpodobností $\mathbf{p}_{(1)}, \dots, \mathbf{p}_{(K)}$ jsou totožné.

Tedy vidíme, že hypotéza homogenity multinomických rozdělení (8.10) je ekvivalentní nezávislosti X a Z (8.8) ve sdruženě multinomickém modelu (8.7). Odtud se dá vyvodit, že testová statistika χ^2 daná vzorcem (8.6) je také vhodná statistika pro testování nulové hypotézy (8.10) ve sloupcově multinomickém modelu (8.9). Dále se dá ukázat, že za vhodných předpokladů na sloupcové rozsahy výběrů (n_{+1}, \dots, n_{+K}) platí za H_0 stejně jako výše $\chi^2 \stackrel{\text{as.}}{\sim} \chi_{(J-1)(K-1)}^2$. Tedy test hypotézy *homogenity multinomických rozdělení* (8.10) ve sloupcově multinomickém modelu můžeme provést úplně stejným způsobem jako test nezávislosti ve sdruženě multinomickém modelu.

Poznámka. Výše uvedené by se dalo také shrnout následovně. Testová statistika (8.6) se dá použít jak pro testování hypotézy nezávislosti (8.8) v modelu (8.7), tak pro testování hypotézy homogenity multinomických rozdělení (8.10) v modelu (8.9). To, který model a hypotézu si vybereme, závisí na interpretovatelnosti modelu a nulové hypotézy pro naše konkrétní data.

Za povšimnutí také stojí, že v případě sloupcově multinomického modelu a hypotézy (8.10) se vlastně jedná o K -výběrový test* na shodu parametrů K multinomických rozdělení).

8.2.1. KONTINGENČNÍ TABULKY 2×2

Zabývejme se nyní speciálním případem $J = 2$ a $K = 2$, kdy obě veličiny mohou nabývat pouze dvou hodnot. Výsledná kontingenční tabulka obsahuje 2×2 četnosti:

	$Z = 1$	$Z = 2$	Σ
$X = 1$	n_{11}	n_{12}	n_{1+}
$X = 2$	n_{21}	n_{22}	n_{2+}
Σ	n_{+1}	n_{+2}	N

	$Z = 1$	$Z = 2$	Σ
$X = 1$	p_{11}	p_{12}	p_{1+}
$X = 2$	p_{21}	p_{22}	p_{2+}
Σ	p_{+1}	p_{+2}	1

Testová statistika je

$$\chi^2 = \sum_{j=1}^2 \sum_{k=1}^2 \frac{\left(n_{jk} - \frac{n_{j+}n_{+k}}{N} \right)^2}{\frac{n_{j+}n_{+k}}{N}}. \quad (8.11)$$

Za platnosti hypotézy nezávislosti má asymptoticky rozdělení χ_1^2 . Hypotézu nezávislosti zamítneme, pokud $\chi^2 \geq \chi_1^2(1 - \alpha)$.

χ^2 -TEST JAKO TEST HOMOGENITY DVOU BINOMICKÝCH ROZDĚLENÍ

Představme si, že veličina Z určuje číslo výběru: máme jeden výběr hodnot náhodné veličiny X z jedinců splňujících $Z = 1$ a druhý výběr náhodné veličiny X z jedinců splňujících $Z = 2$. V prvním výběru bylo n_{11} hodnot $X = 1$ (úspěch) a n_{21} hodnot $X = 2$ (neúspěch), celkem n_{+1} pozorování. Pravděpodobnost úspěchu v 1.výběru je $p_{1(1)} = p_{11}/p_{+1}$. V druhém výběru bylo n_{12} hodnot $X = 1$ (úspěch) a n_{22} hodnot $X = 2$ (neúspěch), celkem n_{+2} pozorování. Pravděpodobnost úspěchu v 2.výběru je $p_{1(2)} = p_{12}/p_{+2}$.

Z předchozího víme, že na χ^2 -test lze nahlížet jako na test shody parametrů $p_{1(1)}$ a $p_{1(2)}$ dvou nezávislých binomických rozdělení $\text{Bi}(n_{+1}, p_{1(1)})$ a $\text{Bi}(n_{+2}, p_{1(2)})$. Tuto situaci jsme vlastně řešili v kapitole 7.2.

Značení zavedené v kapitole 7.2 můžeme snadno převést na značení používané nyní a naopak. Naše kontingenční tabulka přepsaná do značení z kapitoly 7.2 vypadá takto:

* K -výběrově testy pro kvantitativní data budeme uvažovat v kapitole 9.

	$Z = 1$	$Z = 2$	Σ
$X = 1$	X_1	X_2	$X_1 + X_2$
$X = 2$	$n - X_1$	$m - X_2$	$n + m - X_1 - X_2$
Σ	n	m	$n + m$

Rozdíl proti situaci v kapitole 7.2 spočívá v tom, že tam byly oba výběry nezávislé, zatímco nyní uvažujeme jeden výběr z multinomického rozdělení se čtyřmi možnými hodnotami. Tehdy byly rozsahy obou výběrů n, m pevné, nyní jsou to binomické náhodné veličiny a pouze celkový počet pozorování $N = n + m$ je pevný. Znovu jsme narazili na dvě různé formulace dvouvýběrového problému, podobně jako je diskutováno na začátku kapitoly 6 o dvouvýběrových testech pro kvantitativní data. Stejně jako v kapitole 6 i nyní nezáleží z hlediska provedení testu na tom, kterou formulaci používáme a jakým způsobem byla kontingenční tabulka vytvořena. Všechny studované metody platí pro obě dvě formulace.

Kapitola 7.2 vysvětluje, jak porovnat riziko události $[X = 1]$ pro různé hodnoty Z . Můžeme použít tři způsoby porovnání:

- *rozdíl pravděpodobností* $d_X = p_{1(1)} - p_{1(2)}$ odhadneme pomocí $\hat{d} = \frac{n_{11}}{n_{+1}} - \frac{n_{12}}{n_{+2}}$;
- *podíl pravděpodobností* $r_X = p_{1(1)}/p_{1(2)}$ odhadneme pomocí $\hat{r} = \frac{n_{11}n_{+2}}{n_{12}n_{+1}}$;
- *poměr šancí* $o_X = \frac{p_{1(1)}(1-p_{1(2)})}{p_{1(2)}(1-p_{1(1)})}$ odhadneme pomocí $\hat{o} = \frac{n_{11}n_{22}}{n_{12}n_{21}}$ (proto se poměru šancí někdy říká *křížový poměr**).

Metody pro testování těchto parametrů a konstrukci intervalů spolehlivosti jsou uvedeny v kapitole 7.2.

Všimněme si, že nezávislost náhodných veličin X a Z je ekvivalentní jednomu ze vztahů

$$d_X = 0, \quad r_X = 1, \quad o_X = 1.$$

Test na nulovost rozdílu rizik nebo jednotkovost relativního rizika či poměru šancí je v této situaci zároveň testem nezávislosti X a Z .

Poznámka. Dá se ukázat, že v tomto případě pro testovou statistiku χ^2 -testu nezávislosti (8.11) platí

$$\chi^2 = T_d^2,$$

kde T_d je testová statistika pro rozdíl pravděpodobností z (7.5).

8.2.2. KONTINGENČNÍ TABULKY $2 \times K$

Nyní budeme uvažovat speciální případ $J = 2$ a $K \geq 2$. Kontingenční tabulka obsahuje $2 \times K$ četností:

* Angl. *cross ratio*

	$Z = 1$	$Z = 2$	\dots	$Z = K$	Σ
$X = 1$	n_{11}	n_{12}	\dots	n_{1K}	n_{1+}
$X = 2$	n_{21}	n_{22}	\dots	n_{2K}	n_{2+}
Σ	n_{+1}	n_{+2}	\dots	n_{+K}	N

	$Z = 1$	$Z = 2$	\dots	$Z = K$	Σ
$X = 1$	p_{11}	p_{12}	\dots	p_{1K}	p_{1+}
$X = 2$	p_{21}	p_{22}	\dots	p_{2K}	p_{2+}
Σ	p_{+1}	p_{+2}	\dots	p_{+K}	N

Toto je zobecnění situace řešené v kapitole 7.2. Můžeme si ji představit i tak, že máme (po sloupcích) K výběrů z binomického rozdělení s potenciálně různými pravděpodobnostmi úspěchu p_{1k}/p_{+k} nebo máme (po řádcích) dva výběry z multinomického rozdělení s potenciálně různými vektory pravděpodobností

$$\left(\frac{p_{11}}{p_{1+}}, \frac{p_{12}}{p_{1+}}, \dots, \frac{p_{1K}}{p_{1+}} \right)^T \quad \text{a} \quad \left(\frac{p_{21}}{p_{2+}}, \frac{p_{22}}{p_{2+}}, \dots, \frac{p_{2K}}{p_{2+}} \right)^T.$$

TESTOVÁNÍ NEZÁVISLOSTI χ^2 TESTEM

X a Z jsou nezávislé, právě když $p_{1(1)} = p_{1(2)} = \dots = p_{1(K)}$. To vyžaduje, aby pro kterékoli dvě skupiny $Z = k_1$ a $Z = k_2$ byl rozdíl rizik 0 nebo relativní riziko či poměr šancí 1. Zatímco zobecnit testování pomocí rozdílů rizik, jednotkovosti relativního rizika či poměrů šancí na tento případ by vyžadovalo další práci, χ^2 test nezávislosti lze zobecnit snadno.

Pokud platí hypotéza, že X a Z jsou nezávislé náhodné veličiny, pravděpodobnosti $\mathbf{p} = (p_{11}, p_{21}, \dots, p_{1K}, p_{2K})^T$ specifikující multinomické rozdělení vektoru \mathbf{n} jsou funkcemi p_{1+} a $p_{+1}, \dots, p_{+(K-1)}$, celkem K parametrů. Máme tedy že testová statistika je

$$\chi^2 = \sum_{j=1}^2 \sum_{k=1}^K \frac{\left(n_{jk} - \frac{n_{j+} n_{+k}}{N} \right)^2}{\frac{n_{j+} n_{+k}}{N}}.$$

Za platnosti hypotézy nezávislosti má asymptoticky rozdělení χ_{2K-K-1}^2 , tj. χ_{K-1}^2 . Hypotézu nezávislosti zamítneme, pokud $\chi^2 \geq \chi_{K-1}^2(1 - \alpha)$.

Podobně jako jsme se v kapitole 8.2.1 dívali na χ^2 -test nezávislosti jako na test homogenity dvou binomických rozdělení, tak se nyní můžeme na test nezávislosti dívat jako na test hypotézy, že K výběrů z binomického rozdělení má stejné pravděpodobnosti úspěchu (jde tedy o K -výběrový test na binomické rozdělení).

Alternativně pak můžeme na test nahlížet také jako na test, že dva výběry z multinomického rozdělení mají stejné vektory pravděpodobností (jde tedy o dvouvýběrový test na multinomické rozdělení).

Příklad. Předpokládejme, že máme data o nejvyšším dosaženém vzdělání (ZŠ, SŠ a VŠ) a o tom, zda dotyčný je kuřák či nekuřák. Zajímá nás, zda kouření souvisí se vzděláním.

Nulovou hypotézu, že kouření nesouvisí se vzděláním lze nahlížet dvěma ekvivalentními způsoby:

- pro každou ze tří skupin dle dosaženého vzdělání je pravděpodobnost kouření stejná (tj. porovnáваме tři binomická rozdělení);
- složení dle vzdělání se mezi kuřáky a nekuřáky neliší (tj. porovnáваме dvě multinomická rozdělení).

Poznámka. Uvažujme, že pozorujeme K nezávislých náhodných veličin X_1, \dots, X_K , kde $X_k \sim \text{Bi}(n_k, p_k)$ pro každé $k \in \{1, \dots, K\}$. Chceme otestovat hypotézy

$$H_0 : p_1 = \dots = p_K, \quad H_1 : \exists_{k \neq j} p_k \neq p_j.$$

Pro tuto situaci se v literatuře často doporučuje zamítnout nulovou hypotézu, pokud

$$Q \geq \chi_{K-1}^2(1 - \alpha), \quad \text{kde } Q = \frac{1}{\bar{p}(1 - \bar{p})} \sum_{k=1}^K n_k (\hat{p}_k - \bar{p})^2,$$

přičemž $\hat{p}_k = \frac{X_k}{n_k}$, $\bar{p} = \frac{1}{N} \sum_{k=1}^K X_k$ a $N = \sum_{k=1}^K n_k$.

V tomto případě se dá ukázat, že

$$Q = \chi^2,$$

kde χ^2 je testová statistika χ^2 -testu nezávislosti spočteného z následující kontingenční tabulky:

	$Z = 1$	$Z = 2$...	$Z = K$	Σ
$X = 1$	X_1	X_2	...	X_K	$\sum_{k=1}^K X_k$
$X = 2$	$n_1 - X_1$	$n_2 - X_2$...	$n_K - X_K$	$N - \sum_{k=1}^K X_k$
Σ	n_1	n_2	...	n_K	N

Tudíž přístup založený na statistice Q je totožný s χ^2 -testem nezávislosti.

Přípravné příklady ke zkoušce.

Vaše řešení by mělo obsahovat matematický model, hypotézu, testovou statistiku a její přesné (či asymptotické) rozdělení za nulové hypotézy. Dále pak kritický obor nebo vzorec pro výpočet p -hodnoty. Mělo by být také řečeno, zda je daný test přesný nebo asymptotický.

1. Terč je rozdělen do 4 segmentů. Do j -tého segmentu padlo n_j střel, $j = 1, \dots, 4$.
 - (a) Navrhněte test hypotézy, že pravděpodobnost zasažení prvního a druhého segmentu je stejná.
 - (b) Navrhněte test hypotézy, že pravděpodobnost zasažení prvního segmentu je dvakrát větší než pravděpodobnost zasažení čtvrtého segmentu.
 - (c) Sestavte interval spolehlivosti pro podíl pravděpodobností zasažení prvního a čtvrtého segmentu.
2. Ve velkém obchodním domě jsou vedle sebe 3 výtahy. Vedení obchodního domu má údaje o tom, kolikrát byl který výtah během posledního týdne využitý. Jak byste otestovali domněnku, že zákazníci nepreferují žádný z výtahů.
3. Čtyři fakulty (MFF, PřF, FSV a FF) se dohodli, že se pokusí porovnat své studenty z hlediska toho, zda jsou praváci nebo leváci. Každá fakulta zjistila tuto skutečnost u 100 náhodně vybraných studentů. Navrhněte vhodný test hypotézy, že mezi fakultami není rozdíl, co se týká složení praváků a leváků.

9. K -VÝBĚROVÝ PROBLÉM PRO KVANTITATIVNÍ DATA

Dvouvýběrové testy ověřují, jestli se dvě skupiny nezávislých pozorování liší v nějaké charakteristice, nejčastěji ve střední hodnotě. Jak ale porovnat více skupin najednou? Pro kategoriální data (binomické či multinomické rozdělení) jsme problém porovnání několika skupin řešili v minulé kapitole. Nyní budeme studovat tento problém u kvantitativních náhodných veličin.

Máme $K \geq 2$ nezávislých náhodných výběrů (skupin)

$$\begin{aligned} & Y_{11}, \dots, Y_{1n_1} \text{ z rozdělení } F_1, \\ & Y_{21}, \dots, Y_{2n_2} \text{ z rozdělení } F_2, \\ & \vdots \\ \text{a } & Y_{K1}, \dots, Y_{Kn_K} \text{ z rozdělení } F_K. \end{aligned}$$

Pozorování označujeme Y_{ki} , kde k je číslo výběru jdoucí od 1 do K a i je index pozorování v rámci daného výběru běžící od 1 do n_k , kde n_k je rozsah k -tého výběru. Označme $N = \sum_{k=1}^K n_k$ a $\mathbf{n} = (n_1, \dots, n_K)^\top$. Platí $\mathbf{1}_K^\top \mathbf{n} = \sum_{k=1}^K n_k = N$.

K -výběrový problém testuje hypotézu *nulového rozdílu*

$$H_0 : F_1(x) = F_2(x) = \dots = F_K(x), \quad \forall x \in \mathbb{R}, \quad (9.1)$$

proti alternativě, že alespoň mezi dvěma skupinami je nějaký rozdíl, tj.

$$H_1 : \exists k \neq l \exists x \in \mathbb{R} : F_k(x) \neq F_l(x).$$

9.1. ANALÝZA ROZPTYLU (JEDNODUCHÉ TŘÍDĚNÍ)

Budeme předpokládat platnost modelu. Všimněme si, že ten požaduje, aby všechna rozdělení F_1, \dots, F_K měla **shodný rozptyl**.

Podobně jako v případě jednovýběrového a dvouvýběrového t -testu za předpokladu shody rozptylů (viz kapitoly 5.3 a 6.3) bude níže uvedený test přesný za předpokladu normality pozorování a asymptotický bez tohoto předpokladu.

Model:

$$\mathcal{F}_n = \{F_k = N(\mu_k, \sigma^2), \mu_k \in \mathbb{R}, k \in \{1, \dots, K\}, \sigma^2 > 0\} \quad (9.2)$$

nebo

$$\mathcal{F}_{as} = \{F_k \in \mathcal{L}_+^2, k \in \{1, \dots, K\}, \text{přičemž } \text{var}(Y_{11}) = \text{var}(Y_{21}) = \text{var}(Y_{K1}) \stackrel{\text{df}}{=} \sigma^2\}.$$

Všimněme si, že v normálním modelu \mathcal{F}_n se jednotlivé skupiny mohou navzájem lišit pouze střední hodnotou.

Označme μ_k střední hodnotu k -té skupiny, tj. $\mu_k = E Y_{ki}$. Budeme se zabývat otázkou, zdali všechny skupiny mají stejnou střední hodnotu.

Testované parametry: Střední hodnoty $\mu_k = E Y_{ki}$.

Hypotéza a alternativa:

$$H_0 : \mu_1 = \dots = \mu_K, \quad H_1 : \exists_{k \neq l} \mu_k \neq \mu_l.$$

Značení. Nechť

- $Y_{k+} \stackrel{\text{df}}{=} \sum_{i=1}^{n_k} Y_{ki}$ a $\bar{Y}_{k+} \stackrel{\text{df}}{=} n_k^{-1} \sum_{i=1}^{n_k} Y_{ki}$ jsou součty a průměry jednotlivých skupin,
- $Y_{++} \stackrel{\text{df}}{=} \sum_{k=1}^K \sum_{i=1}^{n_k} Y_{ki}$ je celkový součet a $\bar{Y}_{++} \stackrel{\text{df}}{=} N^{-1} \sum_{k=1}^K \sum_{i=1}^{n_k} Y_{ki}$ je celkový průměr.

Všimněte si, že \bar{Y}_{++} je vážený průměr skupinových průměrů \bar{Y}_{k+} s vahami n_k , tj.

$$\bar{Y}_{++} = \frac{\sum_{k=1}^K n_k \bar{Y}_{k+}}{\sum_{k=1}^K n_k}.$$

Označme dále pozorování ve skupinách $\mathbf{Y}_k = (Y_{k1}, \dots, Y_{kn_k})^\top$, $k \in \{1, \dots, K\}$ a všechna pozorování $\mathbf{Y} = (\mathbf{Y}_1^\top, \dots, \mathbf{Y}_K^\top)^\top$.

Náš přístup bude založen na několika druzích součtů čtverců, které zavádí následující definice.

Definice 9.1 Součty čtverců v analýze rozptylu:

- $SS_C \stackrel{\text{df}}{=} \sum_{k=1}^K \sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{++})^2$ nazýváme *celkový součet čtverců*^{*},
- $SS_A \stackrel{\text{df}}{=} \sum_{k=1}^K n_k (\bar{Y}_{k+} - \bar{Y}_{++})^2$ nazýváme *součet čtverců skupin*[†],
- $SS_e \stackrel{\text{df}}{=} \sum_{k=1}^K \sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{k+})^2$ nazýváme *residuální součet čtverců*[‡].

Věta 9.1 Platí

$$SS_C = SS_A + SS_e.$$

Důkaz.

$$\begin{aligned} SS_C &= \sum_{k=1}^K \sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{++})^2 = \sum_{k=1}^K \sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{k+} + \bar{Y}_{k+} - \bar{Y}_{++})^2 \\ &= \sum_{k=1}^K \sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{k+})^2 + \sum_{k=1}^K \sum_{i=1}^{n_k} (\bar{Y}_{k+} - \bar{Y}_{++})^2 + 2 \sum_{k=1}^K \sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{k+})(\bar{Y}_{k+} - \bar{Y}_{++}) \\ &= SS_e + \sum_{k=1}^K n_k (\bar{Y}_{k+} - \bar{Y}_{++})^2 + 2 \sum_{k=1}^K (\bar{Y}_{k+} - \bar{Y}_{++}) \sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{k+}) \end{aligned}$$

* Angl. *total sum of squares* † Angl. *between group sum of squares* ‡ Angl. *residual sum of squares, error sum of squares*

$$= SS_e + SS_A + 0,$$

kde jsme využili toho, že

$$\sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{k+}) = Y_{k+} - n_k \bar{Y}_{k+} = 0, \quad \text{pro } k \in \{1, \dots, K\}.$$

□

Poznámka. SS_C měří celkovou variabilitu dat. Tu můžeme rozložit na variabilitu mezi jednotlivými skupinami vyjadřující jejich vzájemnou odlišnost (SS_A) a variabilitu uvnitř jednotlivých skupin SS_e .

Jelikož \bar{Y}_{k+} je odhadem μ_k a \bar{Y}_{++} je odhadem celkové střední hodnoty (za H_0), bude za platnosti hypotézy SS_A malé vzhledem k SS_e . Pokud je SS_A velké vzhledem k SS_e , znamená to, že se průměry jednotlivých skupin od sebe příliš liší a hypotézu o rovnosti středních hodnot bychom měli zamítnout.

Testová statistika tedy bude porovnávat *variabilitu výběrových průměrů* (SS_A) a *variabilitu uvnitř jednotlivých skupin* (SS_e). V následujícím nejdříve prozkoumáme vlastnosti statistik SS_e a SS_A .

Lemma 9.2 Nechť platí model \mathcal{F}_{as} .

(i) Potom

$$E SS_e = (N - K) \sigma^2.$$

(ii) Pokud dokonce platí model \mathcal{F}_n , tak $\frac{SS_e}{\sigma^2} \sim \chi_{N-K}^2$.

Důkaz. Část (i) Všimněme si, že

$$SS_e = \sum_{k=1}^K \sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{k+})^2 = \sum_{k=1}^K (n_k - 1) S_k^2, \quad (9.3)$$

kde $S_k^2 = \frac{1}{n_k - 1} \sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{k+})^2$ je výběrový rozptyl v k -té skupině. Dle věty 2.6(ii) je S_k^2 nestranným odhadem rozptylu σ^2 . Tudíž

$$E SS_e = \sum_{k=1}^K (n_k - 1) \sigma^2 = (N - K) \sigma^2.$$

Část (ii) Pomocí (9.3) můžeme psát

$$\frac{SS_e}{\sigma^2} = \sum_{k=1}^K \frac{(n_k - 1) S_k^2}{\sigma^2}.$$

Pro $k \in \{1, \dots, K\}$ mají náhodné veličiny $\frac{(n_k - 1) S_k^2}{\sigma^2}$ dle věty 2.8(i) χ^2 -rozdělení o $n_k - 1$ stupních volnosti. Navíc jsou tyto náhodné veličiny nezávislé. Tudíž náhodná veličina $\frac{SS_e}{\sigma^2}$ má χ^2 -rozdělení o $\sum_{k=1}^K (n_k - 1) = N - K$ stupních volnosti.

□

Následující lemma shrnuje vlastnosti statistiky SS_A . Ještě než jej vyslovíme, tak si označme „průměrnou střední hodnotu“

$$\bar{\mu} = E \bar{Y}_{++} = \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} E Y_{ki} = \frac{1}{N} \sum_{k=1}^K n_k \mu_k.$$

Lemma 9.3 Nechť platí model \mathcal{F}_{as} .

(i) Potom

$$E SS_A = \sum_{k=1}^K n_k (\mu_k - \bar{\mu})^2 + (K - 1) \sigma^2.$$

(ii) Pokud navíc platí model \mathcal{F}_n , tak SS_A a SS_e jsou nezávislé.

(iii) Pokud navíc platí model \mathcal{F}_n a ještě také **hypotéza** H_0 , pak $\frac{SS_A}{\sigma^2} \sim \chi_{K-1}^2$.

Důkaz. V důkazu budeme využívat toho, že dle věty 9.1

$$SS_A = SS_C - SS_e. \quad (9.4)$$

Část (i). Spočtěme nejdříve $E SS_C$. Podobně jako ve větě 2.4(ii) můžeme psát

$$SS_C = \sum_{k=1}^K \sum_{i=1}^{n_k} (Y_{ki} - \bar{Y}_{++})^2 = \mathbf{Y}^\top \mathbb{A}_C \mathbf{Y}, \quad \text{kde} \quad \mathbb{A}_C = \mathbb{1}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top. \quad (9.5)$$

Tedy s využitím lemmatu 2.5

$$\begin{aligned} E SS_C &= E \mathbf{Y}^\top \mathbb{A}_C E \mathbf{Y} + \text{tr} \left(\mathbb{A}_C \text{var}(\mathbf{Y}) \right) = \sum_{k=1}^K \sum_{i=1}^{n_k} (\mu_k - \bar{\mu})^2 + \sigma^2 \text{tr}(\mathbb{A}_C) \\ &= \sum_{k=1}^K n_k (\mu_k - \bar{\mu})^2 + \sigma^2 (N - 1). \end{aligned}$$

Dále z lemmatu 9.2 víme, že $E SS_e = (N - K) \sigma^2$. Tedy s využitím (9.4)

$$E SS_A = E SS_C - E SS_e = \sum_{k=1}^K n_k (\mu_k - \bar{\mu})^2 + \sigma^2 (K - 1).$$

Část (ii). Všimněme si, že $\mathbf{Y} \sim N_N(\cdot, \sigma^2 \mathbb{1}_N)$. Jelikož budeme směřovat k využití lemmatu 2.7(ii), tak si potřebujeme nejdříve vyjádřit SS_e a SS_A jako kvadratické formy všech pozorování \mathbf{Y} .

Nejdříve si všimněme, že s využitím (9.3) podobně jako v (9.5) dostáváme

$$SS_e = \sum_{k=1}^K (n_k - 1) S_k^2 = \sum_{k=1}^K \mathbf{Y}_k^\top \left(\mathbb{1}_{n_k} - \frac{1}{n_k} \mathbf{1}_{n_k} \mathbf{1}_{n_k}^\top \right) \mathbf{Y}_k = \mathbf{Y}^\top (\mathbb{1}_N - \mathbb{B}) \mathbf{Y}, \quad (9.6)$$

kde

$$\mathbb{B} = \begin{pmatrix} \frac{1}{n_1} \mathbf{1}_{n_1} \mathbf{1}_{n_1}^\top & \mathbb{0}_{n_1 \times n_2} & \cdots & \mathbb{0}_{n_1 \times n_K} \\ \mathbb{0}_{n_2 \times n_1} & \frac{1}{n_2} \mathbf{1}_{n_2} \mathbf{1}_{n_2}^\top & \cdots & \mathbb{0}_{n_2 \times n_K} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{0}_{n_K \times n_1} & \mathbb{0}_{n_K \times n_2} & \cdots & \frac{1}{n_K} \mathbf{1}_{n_K} \mathbf{1}_{n_K}^\top \end{pmatrix}.$$

Dále pro SS_A s využitím (9.4), (9.5) a (9.6) dostáváme

$$\begin{aligned} SS_A &= SS_C - SS_e = \mathbf{Y}^\top (\mathbb{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top) \mathbf{Y} - \mathbf{Y}^\top (\mathbb{I}_N - \mathbb{B}) \mathbf{Y} \\ &= \mathbf{Y}^\top (\mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top) \mathbf{Y}. \end{aligned} \quad (9.7)$$

Jelikož $\Sigma = \text{var}(\mathbf{Y}) = \sigma^2 \mathbb{I}_N$, tak nyní dle lemmatu 2.7(ii) stačí ověřit, že součin matic $(\mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top)(\mathbb{I}_N - \mathbb{B})$ je nulová matice. Počítejme, tedy

$$\begin{aligned} (\mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top)(\mathbb{I}_N - \mathbb{B}) &= \mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top - \mathbb{B}\mathbb{B} + \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \mathbb{B} \\ &= \mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top - \mathbb{B} + \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top = \mathbb{0}_{N \times N}, \end{aligned}$$

kde jsme využili, toho že

$$\mathbb{B}\mathbb{B} = \mathbb{B} \quad \text{a} \quad \mathbf{1}_N^\top \mathbb{B} = \mathbf{1}_N^\top. \quad (9.8)$$

Část (iii). Nejdříve si všimněme, že statistika SS_A je invariantní vůči posunutí, tj. hodnota SS_A se nezmění, pokud ji budeme počítat z „posunutých“ dat $\tilde{\mathbf{Y}} = \mathbf{Y} - c \mathbf{1}_N$, ať je již $c \in \mathbb{R}$ jakékoliv.

Tedy s využitím (9.7) máme

$$\frac{SS_A}{\sigma^2} = \left(\frac{\mathbf{Y} - \mu \mathbf{1}_N}{\sigma} \right)^\top \left(\mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \right) \left(\frac{\mathbf{Y} - \mu \mathbf{1}_N}{\sigma} \right),$$

kde μ je společná hodnota parametrů μ_1, \dots, μ_K za hypotézy.

Nyní $\frac{\mathbf{Y} - \mu \mathbf{1}_N}{\sigma} \sim N_N(\mathbf{0}, \mathbb{I}_N)$. Jsme tedy v situaci lemmatu A.4 s $\mathbb{A} = \mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top$ a $\Sigma = \mathbb{I}_N$. Zbývá ověřit, že matice $\mathbb{A}\Sigma = \mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top$ je idempotentní. Počítejme, tedy

$$\begin{aligned} \left(\mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \right) \left(\mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \right) &= \mathbb{B}\mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \mathbb{B} - \mathbb{B} \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top + \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \\ &= \mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top + \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top = \mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top, \end{aligned}$$

kde jsme využili (9.8), symetrie matice \mathbb{B} a toho, že $\frac{1}{N} \mathbf{1}_N^\top \mathbf{1}_N = 1$. Tedy $\frac{SS_A}{\sigma^2}$ má dle lemmatu A.4 χ^2 -rozdělení s počtem stupňů volnosti

$$\text{tr} \left(\mathbb{B} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \right) = K - 1.$$

□

Všimněme, si že dle lemmatu 9.2(i) je statistika $\frac{SS_e}{N-K}$ nestranný odhad σ^2 . Naproti tomu dle lemmatu 9.3(i) je $\frac{SS_A}{K-1}$ nestranným odhadem σ^2 **pouze za nulové hypotézy**, kdežto za alternativy platí $E \frac{SS_A}{K-1} > \sigma^2$. To nás přivádí k následujícímu testu.

Testová statistika:

$$F_A = \frac{SS_A/(K-1)}{SS_e/(N-K)}$$

Hypotézu budeme zamítat pro **příliš velké** hodnoty F_A .

Věta 9.4 Nechť platí model \mathcal{F}_n a navíc platí hypotéza H_0 , potom $F_A \sim F_{K-1, N-K}$.

Důkaz. Statistiku F_A můžeme přepsat do tvaru

$$F_A = \frac{\frac{SS_A}{\sigma^2}/(K-1)}{\frac{SS_e}{\sigma^2}/(N-K)}.$$

Z lemmatu 9.3(ii) a z lemmatu 9.2(iii) plyne, že $\frac{SS_A}{\sigma^2} \sim \chi_{K-1}^2$ a $\frac{SS_e}{\sigma^2} \sim \chi_{N-K}^2$. Z lemmatu 9.3(ii) pak plyne nezávislost náhodných veličin $\frac{SS_A}{\sigma^2}$ a $\frac{SS_e}{\sigma^2}$. Tvrzení věty pak plyne z definice F -rozdělení (viz definice 2.5). \square

Pomocí věty 9.4 a úvah před touto větou dostáváme.

Kritický obor:

$$H_0 \text{ zamítneme} \Leftrightarrow F_A \geq F_{K-1, N-K}(1-\alpha),$$

kde $F_{K-1, N-K}(1-\alpha)$ je $(1-\alpha)$ -tý kvantil F -rozdělení s $K-1$ a $N-K$ stupni volnosti.

P-hodnota: $1-G(s)$, kde s je pozorovaná hodnota testové statistiky F_A a G je distribuční funkce rozdělení $F_{K-1, N-K}$.

Poznámka.

- Výše popsaná metoda se nazývá *analýza rozptylu** kvůli tomu, jakým způsobem je sestavena testová statistika (porovnávají se v ní vlastně dva odhady σ^2). **Účelem analýzy rozptylu však není analyzovat rozptyl.** Samotný test se pak nazývá *F-test analýzy rozptylu*.
- V normálním modelu se shodnými rozptyly (tj. v modelu \mathcal{F}_n) je *F-test analýzy rozptylu* přesný test rovnosti středních hodnot v $K \geq 2$ nezávislých výběrech.
- Bez předpokladu normality, ale stále za předpokladu shodnosti rozptylů (tj. v modelu \mathcal{F}_{as}), se dá ukázat, že *F-test analýzy rozptylu* dodržuje hladinu alespoň asymptoticky.[†] To plyne z toho, že se za nulové hypotézy pro $N \rightarrow \infty$ dokáže, že

$$F_A \xrightarrow{d} \frac{\chi_{K-1}^2}{K-1} \quad \text{a také} \quad F_{K-1, N-K}(1-\alpha) \rightarrow \frac{\chi_{K-1}^2(1-\alpha)}{K-1},$$

příčemž druhá konvergence plyne přímo z definice F -rozdělení.

* Angl. *analysis of variance, ANOVA* † V literatuře se toto tvrzení zpravidla odbývá tím, že je to známo. Důkaz lze nalézt například v Kika (2015).

Poznámka. Výsledky analýzy rozptylu se tradičně uvádějí formou tabulky.

Zdroj měnlivosti	Součet čtverců	Stupňů volnosti	Podíl	F
Skupina	SS_A	$K - 1$	$\frac{SS_A}{K-1}$	$\frac{SS_A}{K-1} / \frac{SS_e}{N-K}$
Residuální	SS_e	$N - K$	$\frac{SS_e}{N-K}$	
Celkový	SS_C	$N - 1$		

Tvrzení 9.5 Pokud $K = 2$, pak platí

$$F_A = T_{n_1, n_2}^2,$$

kde F_A je testová statistika analýzy rozptylu a T_{n_1, n_2}^2 je čtverec testové statistiky dvou-výběrového t -testu pro test shody středních hodnot za předpokladu shody rozptylů (viz kapitola 6.3).

Důkaz. S využitím (9.3) lze čítelel testové statistiky F_A psát jako

$$\frac{SS_e}{N - K} = \frac{1}{n_1 + n_2 - 2} ((n_1 - 1)S_1^2 + (n_2 - 1)S_2^2) = S_{n_1, n_2}^2, \quad (9.9)$$

kde S_k^2 je výběrový rozptyl v k -té skupině.

Pro úpravu čítelel testové statistiky F_A si nejdříve všimněme, že

$$\bar{Y}_{1+} - \bar{Y}_{++} = \bar{Y}_{1+} - \frac{n_1 \bar{Y}_{1+} + n_2 \bar{Y}_{2+}}{n_1 + n_2} = \frac{n_2 (\bar{Y}_{1+} - \bar{Y}_{2+})}{n_1 + n_2}.$$

A podobně

$$\bar{Y}_{2+} - \bar{Y}_{++} = \frac{n_1 (\bar{Y}_{2+} - \bar{Y}_{1+})}{n_1 + n_2}.$$

Tedy

$$\begin{aligned} \frac{SS_A}{K - 1} &= n_1 (\bar{Y}_{1+} - \bar{Y}_{++})^2 + n_2 (\bar{Y}_{2+} - \bar{Y}_{++})^2 = \frac{(\bar{Y}_{1+} - \bar{Y}_{2+})^2}{(n_1 + n_2)^2} (n_1 n_2^2 + n_2 n_1^2) \\ &= \frac{n_1 n_2 (\bar{Y}_{1+} - \bar{Y}_{2+})^2}{n_1 + n_2} = \frac{(\bar{Y}_{1+} - \bar{Y}_{2+})^2}{\frac{1}{n_1} + \frac{1}{n_2}}. \end{aligned} \quad (9.10)$$

Pomocí (9.9) a (9.10) pak dostáváme

$$F_A = \frac{SS_A / (K - 1)}{SS_e / (N - K)} = \left(\frac{\bar{Y}_{1+} - \bar{Y}_{2+}}{\sqrt{S_{n_1, n_2}^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \right)^2 = T_{n_1, n_2}^2,$$

což bylo dokázat. □

Pro porovnání dvou skupin je tedy analýza rozptylu ekvivalentní dvouvýběrovému t -testu za předpokladu shody rozptylů (viz kapitola 6.3). Avšak pro $K = 2$ je zpravidla vhodnější použít tento t -test, protože umožňuje testovat i jednostranné hypotézy a lze z něj jednoduše odvodit interval spolehlivosti.

Naopak pro $K > 2$ již nelze mluvit o jednostranných hypotézách a ani již tuto situaci nelze řešit jedním intervalem spolehlivosti.

Poznámka. Analýza rozptylu se dále zobecňuje na vícenásobné třídění. Tato zobecnění se probírají v předmětu *Lineární regrese*. Např. dvojnásobné třídění spočívá v tom, že se pozorování klasifikují do KJ skupin podle dvou kategoriálních veličin s K a J hodnotami. Zajímá nás, zdali některá z obou kategoriálních veličin ovlivňuje střední hodnotu pozorování.

PORUŠENÍ PŘEDPOKLADŮ

Porušení shodnosti rozptylů. V tomto případě F -test analýzy rozptylu nedodrží předepsanou hladinu testu přesně ani asymptoticky. Nicméně publikované simulační studie ukazují, že pokud je počet pozorování ve všech skupinách přibližně stejný, pak skutečná hladina F -testu analýzy rozptylu je blízká předepsané hladině.

Pro případ nestejných rozptylů navrhl již Welch (1951) zobecnění testové statistiky a aproximaci jejího rozdělení. Jde vlastně o zobecnění dvouvýběrového Welchova testu na více výběrů. Testová statistika tohoto testu zohledňuje potenciální různost rozptylů a je dána vzorcem

$$F_w = \frac{\sum_{k=1}^K w_k (\bar{Y}_{k+} - \bar{Y}_w)^2}{K-1} \frac{1}{1 + 2\Lambda(K-2)},$$

kde $w_k = \frac{n_k}{S_k^2}$ je váha, kterou přiřazujeme k -té skupině, $\bar{Y}_w = \frac{\sum_{k=1}^K w_k \bar{Y}_{k+}}{\sum_{k=1}^K w_k}$ je odhad společné střední hodnoty za hypotézy a

$$\Lambda = \frac{\sum_{k=1}^K \frac{1}{n_k-1} \left(1 - \frac{w_k}{\sum_{j=1}^K w_j}\right)^2}{K^2 - 1}$$

je jistá korekce, která je pro velké rozsahy výběrů ve všech skupinách blízká nule.

Dá se ukázat, že za platnosti hypotézy i bez předpokladu shodnosti rozptylů (a také bez předpokladu normality) platí

$$(K-1)F_w \xrightarrow{d} \chi_{K-1}^2$$

kde rozsahy všech výběrů rostou do nekonečna, tj.

$$\min \{n_1, \dots, n_K\} \rightarrow \infty \text{ a zároveň } \frac{n_k}{N} \rightarrow \lambda_k > 0, k \in \{1, \dots, K\}. \quad (9.11)$$

Nicméně podobně jako u Welchova dvouvýběrového t -testu (viz str. 118) se z důvodu opatrnosti doporučuje porovnávat testovou statistiku F_w s kvantily F -rozdělení o $K-1$ a $1/(3\Lambda)$ stupních volnosti.

9.2. MNOHONÁSOBNÁ POROVNÁVÁNÍ

V analýze rozptylu porovnáváme mezi sebou střední hodnoty K skupin. Pokud F -test analýzy rozptylu zamítne hypotézu, že všechny skupiny mají stejnou střední hodnotu, pak usoudíme, že alespoň některé skupiny se od sebe liší ve středních hodnotách. Nevíme ovšem, kolik takových odlišných skupin je, ani které to jsou.

Kdybychom chtěli porovnat střední hodnoty pouze dvou skupin, třeba skupin k a l , použili bychom dvouvýběrový t -test. Mohli bychom pak provést dvouvýběrové testy pro všech $\frac{K(K-1)}{2}$ možných dvojic skupin a otestovat všechny hypotézy $H_0^{kl} : \mu_k = \mu_l$ na hladině α . Potom ale pravděpodobnost, že alespoň jednu hypotézu zamítneme za podmínky, že všechny hypotézy platí, není rovna α , ale je větší.

Problém současného testování více hypotéz se ve statistice často nazývá *problém mnohonásobných porovnávání** nebo *mnohonásobného testování†*.

Obecný problém mnohonásobného testování můžeme formulovat následovně. Máme m hypotéz H_0^1, \dots, H_0^m , které chceme otestovat. Pro test hypotézy H_0^j použijeme testovou statistiku T_j a kritickým oborem C_j zvoleným tak, aby každý test měl hladinu α_0 . Pro každé $j \in \{1, \dots, m\}$ tedy platí

$$P_{H_0^j} [T_j \in C_j] = \alpha_0.$$

Celková pravděpodobnost zamítnutí alespoň jedné hypotézy za předpokladu, že všechny platí, je

$$P_{\bigcap_{j=1}^m H_0^j} \left(\bigcup_{j=1}^m [T_j \in C_j] \right) = \alpha_C.$$

Pochopitelně α_C je větší než α_0 , často výrazně. Naším cílem tedy je pro předepsané α najít takové testy \tilde{T}_j s kritickými obory \tilde{C}_j , aby

$$P_{\bigcap_{j=1}^m H_0^j} \left(\bigcup_{j=1}^m [\tilde{T}_j \in \tilde{C}_j] \right) \leq \alpha.$$

Podobně pro intervaly spolehlivosti. Nechť B_1, \dots, B_m jsou intervaly spolehlivosti pro parametry $\theta_X^{(1)}, \dots, \theta_X^{(m)}$, které splňují

$$P(B_j \ni \theta_X^{(j)}) = 1 - \alpha, \quad j \in \{1, \dots, m\}.$$

kde $1 - \alpha$ je předepsaná pravděpodobnost pokrytí.

Potom typicky

$$P(B_1 \ni \theta_X^{(1)}, \dots, B_m \ni \theta_X^{(m)}) < 1 - \alpha.$$

Naším cílem je sestavit intervaly spolehlivosti $\tilde{B}_1, \dots, \tilde{B}_m$, pro které platí

$$P(\tilde{B}_1 \ni \theta_X^{(1)}, \dots, \tilde{B}_m \ni \theta_X^{(m)}) \geq 1 - \alpha.$$

Intervaly $\tilde{B}_1, \dots, \tilde{B}_m$ pak nazýváme *simultánní intervaly spolehlivosti‡*.

V této kapitole si uvedeme nejprve jeden obecný přístup k tomuto problému a pak speciální metodu pro porovnávání středních hodnot několika nezávislých výběrů.

* Angl. *multiple comparisons* † Angl. *multiple testing* ‡ Angl. *simultaneous confidence intervals*

9.2.1. BONFERRONIHO METODA

Máme danou celkovou hladinu α a chceme zaručit, že $\alpha_C \leq \alpha$. K tomu použijeme následující lemma.

Lemma 9.6 (Booleova nerovnost) Pro libovolné náhodné jevy A_1, \dots, A_m platí

$$P\left(\bigcup_{j=1}^m A_j\right) \leq \sum_{j=1}^m P(A_j).$$

Booleova nerovnost je triviální pro $m = 2$, pro vyšší m se snadno dokáže matematickou indukcí.

Máme tedy

$$\alpha_C = P_{\bigcap_{j=1}^m H_0^j} \left(\bigcup_{j=1}^m [T_j \in C_j] \right) \leq m\alpha_0.$$

Zvolíme-li $\alpha_0 = \alpha/m$, pak musí platit $\alpha_C \leq \alpha$. Chceme-li tedy provést m testů tak, aby celková hladina všech testů (pravděpodobnost zamítnutí alespoň jedné hypotézy za podmínky, že všechny platí) byla nejvýše α , provedeme jednotlivé dílčí testy na hladině α/m .

Podobně, chceme-li sestavit m intervalů spolehlivosti tak, aby pravděpodobnost, že všechny intervaly pokryjí hledané parametry, byla alespoň $1 - \alpha$, stačí stanovit pravděpodobnost pokrytí jednotlivých dílčích intervalů na $1 - \alpha/m$. Tento přístup k mnohonásobnému testování a konstrukci simultánních intervalů spolehlivosti se nazývá *Bonferroniho metoda**

Výhodou Bonferroniho metody je její jednoduchost a universalita. Její nevýhodou je, že úprava hladiny α na α/m je téměř vždy příliš přísná. Bonferroniho metoda tedy dává testy s malou silou a zbytečně široké intervaly spolehlivosti. Tyto nevýhody se snaží napravit různé speciální metody mnohonásobného porovnávání, které byly odvozeny pro konkrétní problémy (viz např. Tukeyho metoda popsaná níže).

Aplikace Bonferroniho metody na mnohonásobná porovnávání v analýze rozptylu vypadá takto: provedeme $\frac{K(K-1)}{2}$ dvouvýběrových t -testů pro všechny možné dvojice skupin a otestujeme všechny hypotézy $H_0^{kl} : \mu_k = \mu_l$ na hladině $\frac{2\alpha}{K(K-1)}$. Pokud je některá z těchto hypotéz zamítnuta, prohlásíme střední hodnoty daných dvou skupin za významně odlišné na celkové hladině α .

Máme-li například $\alpha = 0,05$ a $K = 6$ skupin, provádíme 15 testů rovnosti středních hodnot pro 15 dvojic různých skupin na hladině $0,05/15 \doteq 0,0033$. To je natolik nízká hladina, že může být obtížné najít kterékoli dvě odlišné skupiny, přestože F -test analýzy rozptylu zamítá hypotézu, že všechny střední hodnoty jsou stejné.

Poznámka. Při používání metod, které berou v potaz problematiku mnohonásobného testování, se někdy zavádí tzv. *upravená/korigovaná p-hodnota*[†]. Tato upravená p-hodnota se v případě Bonferroniho metody spočte jednoduše jako

$$\tilde{p}_j = \min \{mp_j, 1\}, \quad j \in \{1, \dots, m\},$$

* Angl. *Bonferroni correction* † Angl. *p-value adjusted for multiple comparison*

kde p_j je klasická (neupravená) p -hodnota j -tého testu.

9.2.2. TUKEYOVA METODA

Tato metoda vychází z normálního (homoskedastického) modelu (9.2) předpokládaném při analýze rozptylu. Za platnosti tohoto modelu pak ve srovnání s Bonferroniho metodou dává testy s vyšší silou a kratší intervaly spolehlivosti.

Mějme nezávislé náhodné veličiny $Z_j \sim N(\mu, \sigma^2)$ pro $j \in \{1, \dots, m\}$. Nechť S^2 je odhad rozptylu σ^2 takový, že S^2 je nezávislé na Z_1, \dots, Z_m a pro nějaké přirozené v platí $\frac{vS^2}{\sigma^2} \sim \chi_v^2$.

Definujme tak řečené *studentisované rozpětí** jako

$$Q = \frac{\max_{j \in \{1, \dots, m\}} Z_j - \min_{j \in \{1, \dots, m\}} Z_j}{S}.$$

Lze ukázat, že náhodná veličina Q má rozdělení závislé pouze na hodnotách m a v . Označme kvantilovou funkci tohoto rozdělení $q_{m,v}(\alpha)$. (Vzorce pro hustotu a kvantilovou funkci studentisovaného rozpětí nebudeme uvádět.)[†]

Studentizovaného rozpětí lze použít k sestrojení simultánních intervalů spolehlivosti pro rozdíly středních hodnot. Tento postup se nazývá *Tukeyova metoda*.[‡]

Věta 9.7 (Tukeyova) Nechť Z_1, \dots, Z_m jsou nezávislé náhodné veličiny s rozdělením $Z_j \sim N(\mu_j, \sigma^2)$. Nechť S^2 je odhad rozptylu σ^2 takový, že S^2 je nezávislé na Z_1, \dots, Z_m a pro nějaké přirozené v platí $\frac{vS^2}{\sigma^2} \sim \chi_v^2$. Pak

$$P\left[Z_k - Z_j - Sq_{m,v}(1-\alpha) \leq \mu_k - \mu_k \leq Z_k - Z_j + Sq_{m,v}(1-\alpha), \forall k \neq j \in \{1, \dots, m\}\right] = 1 - \alpha.$$

Tukeyovu větu lze snadno použít i na testování hypotéz. Hypotézu $H_0^{kj} : \mu_k = \mu_j$ zamítneme, pokud $|Z_k - Z_j| > Sq_{m,v}(1-\alpha)$. Hypotézu $H_0 : \mu_1 = \dots = \mu_m$ zamítneme na celkové hladině α , pokud pro alespoň jednu dvojici $k \neq j$ platí $|Z_k - Z_j| > Sq_{m,v}(1-\alpha)$.

Tukeyovu větu můžeme přímo aplikovat na mnohonásobná porovnávání v analýze rozptylu, pokud rozsah výběru všech skupin je totožný, tj. $n_1 = \dots = n_K \equiv n$. Pak totiž $\bar{Y}_{1+}, \dots, \bar{Y}_{K+}$ jsou nezávislé náhodné veličiny s rozdělením $\bar{Y}_{k+} \sim N(\mu_k, \frac{\sigma^2}{n})$. Za S^2 , tj. za odhad $\frac{\sigma^2}{n}$ vezmeme $\frac{SS_e}{n(N-K)}$. Máme $v = N - K$. Hypotézu $H_0^{kj} : \mu_k = \mu_j$ zamítneme, pokud

$$\left|\bar{Y}_{k+} - \bar{Y}_{j+}\right| \geq \sqrt{\frac{SS_e}{N-K}} \sqrt{\frac{1}{n}} q_{K, N-K}(1-\alpha). \quad (9.12)$$

Pokud rozsahy všech výběrů nejsou stejné, nemůžeme Tukeyovu větu přímo použít, protože nejsou splněny její předpoklady. Lze ale dokázat, že pokud výraz $\sqrt{\frac{1}{n}}$ v (9.12)

* Angl. *studentized range* † Studentisované rozpětí se někdy definuje jako $Q/\sqrt{2}$. Na to je třeba dávat pozor při používání tabelovaných nebo softwarem vypočtených hodnot $q_{m,v}(\alpha)$. Pro kontrolu můžeme porovnat rozdělení Q při $m = 2$ s rozdělením t_k . Pro naši definici jsou tato dvě rozdělení totožná.

‡ Angl. *Tukey method, Tukey's range test, Tukey's HSD (honest significant difference) test.*

nahradíme výrazem $\sqrt{\frac{1}{2n_k} + \frac{1}{2n_j}}$, celková pravděpodobnost zamítnutí některé z platných hypotéz H_0^{kj} nepřekročí α . Tukeyova metoda tedy po této úpravě stále funguje, pouze se stává poněkud konservativní.

9.3. KRUSKALŮV-WALLISŮV TEST

Kruskalův-Wallisův test je zobecněním dvouvýběrového Wilcoxonova testu na porovnání $K \geq 2$ výběrů. I nadále používáme značení zavedené na začátku kapitoly pro K -výběrový problém.

Model: $\mathcal{F} = \{\exists g(\cdot)$ rostoucí funkce $\exists F$ spojitá d.f. $\exists \delta_1, \dots, \delta_K \in \mathbb{R}$:

$$g(X_{k1}) \sim F_k, F_k(x) = F(x - \delta_k) \forall x \in \mathbb{R}, k \in \{1, \dots, K\}$$

Jde o K spojitých rozdělání, které jsou po nějaké vhodné rostoucí transformaci g navzájem posunutá v poloze.

Hypotéza a alternativa:

$$H_0 : \delta_1 = \dots = \delta_K, \quad H_1 : \exists_{k \neq l} \delta_k \neq \delta_l.$$

Poznámka. Pokud platí model \mathcal{F} a hypotéza H_0 , rozdělání ve všech skupinách jsou totožná. Potom platí mezi K skupinami rovnost veškerých charakteristik (tj. platí hypotéza nulového rozdílu (9.1)).

Náznak odvození testové statistiky

Lze ukázat, že testová statistika dvouvýběrového Wilcoxonova testu (6.8) je ekvivalentní čitateli testové statistiky dvouvýběrového t -testu (tj. rozdílu průměrů), pokud do ní místo původních pozorování dosadíme jejich pořadí. Se stejnou logikou můžeme zkusit postupovat jako při konstrukci F -testu analýzy rozptylu, kde hodnoty původních pozorování ve spojeném (sdruženém) výběru $\mathbf{Y} = (Y_{11}, \dots, Y_{Kn_K})^T$ nahradíme jejich pořadími R_{11}, \dots, R_{Kn_K} (tj. R_{ki} je pořadí Y_{ki} mezi všemi pozorováními \mathbf{Y}).

Potom

$$\widetilde{SS}_A = \sum_{k=1}^K n_k (\bar{R}_{k+} - \bar{R}_{++})^2,$$

kde $\bar{R}_{k+} = \frac{1}{n_k} \sum_{i=1}^{n_k} R_{ki}$ je průměrné pořadí v k -té skupině a

$$\bar{R}_{++} = \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} R_{ki} = \frac{N+1}{2}$$

je celkové průměrné pořadí.

Dále se všimněme, že ve standardní analýze rozptylu veličina $\frac{SS_e}{N-K}$ odhaduje neznámý rozptyl $\text{var}(Y_{ki}) = \sigma^2$. V případě pořadí však díky větě 2.16(iii) víme, že za hypotézy $\tilde{\sigma}^2 = \text{var}(R_{ki}) = \frac{N^2-1}{12}$. Jako testová statistika se pak nabízí

$$\tilde{Q} = \frac{\tilde{SS}_A}{\tilde{\sigma}^2} = \frac{12}{(N-1)(N+1)} \sum_{k=1}^K n_k \left(\bar{R}_{k+} - \frac{N+1}{2} \right)^2. \quad (9.13)$$

Dá se ukázat*, že v asymptotickém smyslu platí analogie lemmatu 9.3(iii), tj. za hypotézy a při rostoucím počtu pozorování, viz (9.11),

$$\tilde{Q} \xrightarrow{d} \chi_{K-1}^2.$$

Jak ale ukážeme níže, tak platí $E \tilde{Q} = (K-1) \frac{N}{N-1}$. Protože však limitní χ_{K-1}^2 -rozdělení má střední hodnotu $K-1$, tak pro zlepšení asymptotické aproximace se používá $Q = \frac{N-1}{N} \tilde{Q}$.

Testová statistika:

$$Q = \frac{12}{N(N+1)} \sum_{k=1}^K n_k \left(\bar{R}_{k+} - \frac{N+1}{2} \right)^2.$$

Kritický obor (asymptotický): Jelikož proti nulové hypotézy svědčí velké hodnoty testové statistiky, tak pro asymptotický test dostáváme

$$H_0 \text{ zamítneme} \Leftrightarrow Q \geq \chi_{K-1}^2(1-\alpha).$$

Výše uvedený test se nazývá *Kruskalův-Wallisův test*. Podobně jako u (dvouvýběrového) Wilcoxonova testu se pro malé rozsahy výběru (v případě, že v datech nejsou shody) používají přesné kritické hodnoty, které jsou tabelovány.

Poznámka. Položme $R_{k+} = \sum_{i=1}^{n_k} R_{ki}$. Potom

$$\begin{aligned} \sum_{k=1}^K n_k \left(\bar{R}_{k+} - \frac{N+1}{2} \right)^2 &= \sum_{k=1}^K \frac{1}{n_k} \left(R_{k+} - n_k \frac{N+1}{2} \right)^2 \\ &= \sum_{k=1}^K \frac{1}{n_k} \left(R_{k+}^2 - R_{k+} n_k (N+1) + n_k^2 \frac{(N+1)^2}{4} \right) = \sum_{k=1}^K \frac{R_{k+}^2}{n_k} - \frac{N(N+1)^2}{4}. \end{aligned}$$

Proto se často testová statistika Q uvádí ve výpočetně jednodušší formě (viz např. [Anděl, 2002](#), kapitola 11.3.1)

$$Q = \frac{12}{N(N+1)} \sum_{k=1}^K \frac{R_{k+}^2}{n_k} - 3(N+1). \quad (9.14)$$

* Viz např. [Hájek and Šidák \(1967\)](#).

Poznámka. Využijme nyní vzorce (9.14) pro výpočet střední hodnoty statistiky Q za hypotézy. Za tímto účelem nejdříve s využitím věty 2.16 spočítáme

$$\begin{aligned} E R_{k+}^2 &= \text{var} (R_{k+}) + (E R_{k+})^2 \\ &= \sum_{i=1}^{n_k} \text{var} (R_{ki}) + \sum_{i=1}^{n_k} \sum_{i'=1, i' \neq i}^{n_k} \text{cov} (R_{ki}, R_{ki'}) + \left(\sum_{i=1}^{n_k} \frac{N+1}{2} \right)^2 \\ &= \frac{n_k(N^2-1)}{12} - \frac{n_k(n_k-1)(N+1)}{12} + \frac{n_k^2(N+1)^2}{4}. \end{aligned}$$

Tudíž (za hypotézy)

$$\begin{aligned} E Q &= \frac{12}{N(N+1)} \sum_{k=1}^K \frac{E R_{k+}^2}{n_k} - 3(N+1) \\ &= \frac{12}{N(N+1)} \sum_{k=1}^K \left[\frac{(N^2-1)}{12} - \frac{(n_k-1)(N+1)}{12} + \frac{n_k(N+1)^2}{4} \right] - 3(N+1) \\ &= \frac{K(N-1)}{N} - \frac{(N-K)}{N} + 3(N+1) - 3(N-1) = K-1, \end{aligned}$$

což odpovídá střední hodnotě rozdělení χ_{K-1}^2 a vysvětluje, proč se místo testové statistiky \tilde{Q} dané vzorcem (9.13) používá statistika Q .

PORUŠENÍ PŘEDPOKLADŮ

Shody kvůli zaokrouhlování. Kvůli zaokrouhlování bývají v aplikacích v datech často shody. Testová statistika Q se pak spočte s využitím tzv. *průměrných pořadí*, viz. (2.6). Dá se ukázat, že potom za nulové hypotézy

$$\frac{Q}{1 - kor.} \xrightarrow{(9.11)} \chi_{K-1}^2,$$

kde *kor.* je korekce upravující rozptyl daná předpisem*

$$kor. = \frac{1}{N(N^2-1)} \sum_y (t_y^3 - t_y),$$

přičemž t_y značí počet, kolikrát se mezi všemi hodnotami $Y_{11}, \dots, Y_{K n_K}$ vyskytla hodnota y . Za povšimnutí stojí, že bez této úpravy by byl test (asymptoticky) konzervativní.

Neplatí zobecněný model posunutí. Nejdříve si všimněme, že test (asymptoticky) dodržuje hladinu, pokud jsou všechna pozorování $Y_{11}, \dots, Y_{K n_K}$ nezávislá a stejně rozdělená. Neplatnost modelu posunutí má tedy, podobně jako u dvouvýběrového Wilcoxonova testu (viz kapitola 6.4), dva nepříjemné důsledky pro chování testu za alternativy:

* Viz např. Hollander et al. (2013), str. 205.

1. **interpretační problém** - jestliže neplatí zobecněný model posunutí, pak ze zamítnutí nulové hypotézy lze vyvodit pouze to, že rozdělení v jednotlivých skupinách nejsou shodná. Obecně však nelze vyvozovat, že se liší střední hodnoty, resp. mediány skupin.
2. **síla testu** - podobně jako u Mannovy-Whitneyho formulace dvouvýběrového Wilcoxonova testu (viz str. 126) lze ukázat, že Kruskalův-Wallisův test se zaměřuje na zkoumání, zda pro všechna $k, l \in \{1, \dots, K\}$ platí $P[Y_{k1} < Y_{l1}] = 1/2$. Pokud v zobecněném modelu posunutí je $\delta_k \neq \delta_l$, pak vskutku $P[Y_{k1} < Y_{l1}] \neq 1/2$. Pokud však za alternativy dojde k jiným změnám než pouze v parametrech polohy, pak není zřejmé, jaký to bude mít důsledek na sílu testu.

Přípravné příklady ke zkoušce.

Vaše řešení „praktických úloh“ by mělo obsahovat matematický model, hypotézu, testovou statistiku a její přesné (či asymptotické) rozdělení za nulové hypotézy. Dále pak kritický obor nebo vzorec pro výpočet p -hodnoty. Mělo by být také řečeno, zda je daný test přesný nebo asymptotický.

1. Máme údaje o tělesné výšce 500 dospělých žen a jejich barvě očí, přičemž rozlišujeme hnědou, modrou a zelenou barvu. Navrhněte vhodný test, který by zjistil, zda tělesná výška souvisí s barvou očí.
2. Máme údaje o platech 2 000 pracovníků v oboru IT a regionu (8 možností), ve kterém žijí. Navrhněte vhodný postup, který by při dodržení předepsané celkové hladiny hledal dvojice regionů, pro které lze vyšší platu považovat za různou.

10. KORELAČNÍ ANALÝZA

Uvažujme náhodný výběr

$$\begin{pmatrix} X_1 \\ Y_1 \end{pmatrix}, \dots, \begin{pmatrix} X_n \\ Y_n \end{pmatrix}$$

dvousložkových náhodných vektorů a $n \geq 3$.

10.1. PEARSONŮV KORELAČNÍ KOEFICIENT

Chceme otestovat nezávislost mezi X_i a Y_i , případně sestavit interval spolehlivosti pro korelační koeficient definovaný jako

$$\rho = \rho(X_i, Y_i) = \frac{\text{cov}(X_i, Y_i)}{\sqrt{\text{var } X_i \text{ var } Y_i}}. \quad (10.1)$$

Jeho konsistentním odhadem je výběrový korelační koeficient

$$\begin{aligned} \hat{\rho}_n &= \frac{S_{XY}}{S_X S_Y} = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)(Y_i - \bar{Y}_n)}{\sqrt{\sum_{i=1}^n (X_i - \bar{X}_n)^2 \sum_{i=1}^n (Y_i - \bar{Y}_n)^2}} \\ &= \frac{\sum_{i=1}^n X_i Y_i - n \bar{X}_n \bar{Y}_n}{\sqrt{\left(\sum_{i=1}^n X_i^2 - n \bar{X}_n^2\right) \left(\sum_{i=1}^n Y_i^2 - n \bar{Y}_n^2\right)}} \end{aligned} \quad (10.2)$$

zavedený v definici 3.8.

Poznámka. Připomeňme, že $\rho = 1$ (resp. -1), právě když existují konstanty $a \in \mathbb{R}$ a $b > 0$ (resp. $b < 0$) takové, že s pravděpodobností jedna platí $Y = a + bX$.

Podobně $\hat{\rho}_n = 1$ (resp. -1) právě když existují konstanty $a \in \mathbb{R}$ a $b > 0$ (resp. $b < 0$) takové, že $Y_i = a + bX_i$ pro všechna $i = 1, \dots, n$.

Většina základních tvrzení o korelačním koeficientu jsou odvozena za předpokladu normality.

Model:

$$\mathcal{F}_n = \left\{ \begin{pmatrix} X_i \\ Y_i \end{pmatrix} \sim N_2 \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{pmatrix} \right), \mu_1, \mu_2 \in \mathbb{R}, \sigma_1^2 > 0, \sigma_2^2 > 0, \rho \in \langle -1, 1 \rangle \right\}.$$

Některá tvrzení však platí i v obecnějším modelu

$$\mathcal{F}_{as} = \{ \text{dvourozměrná rozdělení přičemž } \text{var}(X_i) \in (0, \infty), \text{var}(Y_i) \in (0, \infty) \}.$$

10.1.1. TESTOVÁNÍ HYPOTÉZY NEZÁVISLOSTI

Hypotéza a alternativa:

$$H_0 : X_i \text{ a } Y_i \text{ jsou nezávislé} \quad H_1 : X_i \text{ a } Y_i \text{ nejsou nezávislé.} \quad (10.3)$$

Tvrzení 10.1 Nechť platí model \mathcal{F}_n a navíc jsou složky X_i a Y_i **nezávislé** (tj. $\rho = 0$). Potom

$$T_n = \sqrt{n-2} \frac{\hat{\varrho}_n}{\sqrt{1-\hat{\varrho}_n^2}} \sim t_{n-2}.*$$

Všimněme si, že za hypotézy nezávislosti (10.3) je korelační koeficient ρ nulový. Tudíž proti hypotéze nezávislosti svědčí hodnoty výběrového korelačního koeficientu $\hat{\varrho}_n$ příliš vzdálené od nuly. Tedy díky tvrzení 10.1 dostáváme kritický obor:

$$H_0 \text{ zamítneme} \Leftrightarrow |T_n| \geq t_{n-2}(1-\alpha/2). \quad (10.4)$$

Tento test má (za platnosti modelu \mathcal{F}_n) přesně hladinu α .

Poznámka.

- Za platnosti modelu \mathcal{F}_n je test nezávislosti (10.3) ekvivalentní s testem nekorelovanosti, tj.

$$H_0 : \varrho = 0 \text{ proti alternativě } H_1 : \varrho \neq 0.$$

Jelikož závislost v dvourozměrném normálním rozdělení implikuje, že $\rho \neq 0$, tak se dá ukázat, že výše popsaný test je v modelu \mathcal{F}_n konzistentní.

- Dá se ukázat, že pro platnost tvrzení 10.1 stačí, že pouze jedna z veličin má normální rozdělení a druhá je s ní nezávislá a má rozdělení s konečným a nenulovým rozptylem.

Nechť $\begin{pmatrix} X_i \\ Y_i \end{pmatrix}$, $i = 1, \dots, n$ je náhodný výběr z dvourozměrného rozdělení z modelu \mathcal{F}_{as} (tj. nepředpokládáme normalitu) Potom se dá, za předpokladu, že X_i a Y_i jsou **nezávislé**, ukázat

$$T_n = \sqrt{n-2} \frac{\hat{\varrho}_n}{\sqrt{1-\hat{\varrho}_n^2}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1). \quad (10.5)$$

Tj. bez předpokladu normality je test s kritickým oborem (10.4) asymptotickým testem hypotézy nezávislosti (10.3). Tento test bude citlivý proti alternativám (tj. konzistentní pro alternativy), pro které je skutečný korelační koeficient ϱ nenulový, tj. mezi veličinami se dá detekovat lineární závislost. Na druhou stranu test však nebude konzistentní, pokud X_i a Y_i sice nebudou nezávislé, ale jsou nekorelované, tj. $\varrho = 0$. Extrémní příkladem takové situace je, když X_i má symetrické rozdělení kolem nuly a $Y_i = X_i^2$.

Je důležité si uvědomit, že se zde jedná o test nezávislosti nikoliv však o test nekorelovanosti, tj. o test hypotézy, že $H_0 : \varrho = 0$. Existují totiž dvourozměrná rozdělení, pro která je $\varrho = 0$, ale asymptotický výsledek (10.5) neplatí.

* Důkaz, viz věta 6.2 v Anděl (2002).

10.1.2. TESTOVÁNÍ OBECNÉ HODNOTY KORELAČNÍHO KOEFICIENTU A INTERVAL SPOLEHLIVOSTI PRO ρ

Nechť platí model \mathcal{F}_n a $\rho_0 \neq 0$. Všimněme si, že tvrzení 10.1 (ani výsledek (10.5)) nelze rozšířit na testování hypotéz

$$H_0 : \rho = \rho_0, \quad H_1 : \rho \neq \rho_0. \quad (10.6)$$

Na základě těchto tvrzení nelze také sestavit interval spolehlivosti.

I když se dá za platnosti modelu \mathcal{F}_n nalézt (přesné) rozdělení výběrového korelačního koeficientu $\widehat{\rho}_n$, tak toto rozdělení je analyticky natolik složité, že se v praxi nepoužívá. Nicméně se pomocí Δ -metody (tvrzení 1.7) dá ukázat, že

$$\sqrt{n} (\widehat{\rho}_n - \rho) \xrightarrow[n \rightarrow \infty]{d} N(0, (1 - \rho^2)^2). \quad (10.7)$$

Odtud bychom mohli sestavit asymptotický interval spolehlivosti pro parametr ρ jako

$$\left(\widehat{\rho}_n - \frac{u_{1-\alpha/2}(1-\widehat{\rho}_n^2)}{\sqrt{n}}, \widehat{\rho}_n + \frac{u_{1-\alpha/2}(1-\widehat{\rho}_n^2)}{\sqrt{n}} \right).$$

Ukazuje se však, že pro ρ blízké plus nebo minus jedné je lépe použít transformaci stabilizující asymptotický rozptyl

$$\operatorname{arctgh} x = \frac{1}{2} \log \frac{1+x}{1-x},$$

kteřou navrhl R. A. Fisher. Tato transformace se nazývá *Fisherova Z-transformace*.

Tvrzení 10.2 Nechť platí model \mathcal{F}_n . Potom

$$\sqrt{n-3} \left(\operatorname{arctgh} \widehat{\rho}_n - \operatorname{arctgh} \rho \right) \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Důkaz. Důkaz plyne z asymptotické normality korelačního koeficientu (10.7) a Δ -metody (tvrzení 1.7). \square

Poznámka. Použití $\sqrt{n-3}$ místo tradičního \sqrt{n} nemá na asymptotické rozdělení vliv. Ukazuje se však, že normální aproximace je pro $\sqrt{n-3}$ přesnější než pro \sqrt{n} . Důvodem je, že v modelu \mathcal{F}_n se dá odvodit, že

$$\operatorname{var} (\operatorname{arctgh} \widehat{\rho}_n) = \frac{1}{n-3} + o\left(\frac{1}{n}\right).$$

Chceme-li tedy otestovat hypotézu (10.6), tak spočítáme testovou statistiku

$$Z_n = \sqrt{n-3} (\operatorname{arctgh} \widehat{\rho}_n - \operatorname{arctgh} \rho_0)$$

a H_0 zamítneme na hladině α , pokud $|Z_n| \geq u_{1-\alpha/2}$. Z tvrzení 10.2 dostáváme asymptotický interval spolehlivosti pro $\operatorname{arctgh} \rho$ ve tvaru

$$\left(\operatorname{arctgh} \widehat{\rho}_n - \frac{u_{1-\alpha/2}}{\sqrt{n-3}}, \operatorname{arctgh} \widehat{\rho}_n + \frac{u_{1-\alpha/2}}{\sqrt{n-3}} \right).$$

Odtud pak zpětnou transformací pomocí funkce $\operatorname{tgh} x = \frac{\exp(2x)-1}{\exp(2x)+1}$ dostaneme asymptotický interval spolehlivosti pro ρ

$$\left(\operatorname{tgh} \left(\operatorname{arctgh} \widehat{\rho}_n - \frac{u_{1-\alpha/2}}{\sqrt{n-3}} \right), \operatorname{tgh} \left(\operatorname{arctgh} \widehat{\rho}_n + \frac{u_{1-\alpha/2}}{\sqrt{n-3}} \right) \right).$$

Zdůrazněme, že Fisherova Z-transformace spoléhá na **dvourozměrnou normalitu**. Pro jiná rozdělení tvrzení 10.2 obecně neplatí. Pokud nechceme předpokládat, že data pocházejí z dvourozměrného normálního rozdělení, tak je třeba pomocí Δ -metody (tvrzení 1.7) najít asymptotické rozdělení výběrového korelačního koeficientu. Toto asymptotické rozdělení pak má mnohem složitější asymptotický rozptyl než je uvedeno v (10.7). Alternativně lze také využít metodu *bootstrap*, s níž se lze seznámit v předmětu *NMST545 Matematická statistika 4*.

10.2. SPEARMANŮV KORELAČNÍ KOEFICIENT

Spearmanův korelační koeficient vychází z výrazu (10.2), ale dosazuje do něj pořadí namísto původních pozorování. Označme R_i pořadí pozorování X_i v náhodném výběru X_1, \dots, X_n a označme S_i pořadí pozorování Y_i v náhodném výběru Y_1, \dots, Y_n . Pokud X_i je nezávislé na Y_i pro každé i , pak by neměly být závislosti ani mezi pořadími R_i a S_i .

(Výběrový) *Spearmanův korelační koeficient** dostaneme dosazením R_i místo X_i a S_i místo Y_i v (10.2):

$$\widehat{\rho}_n^{(S)} = \frac{\sum_{i=1}^n R_i S_i - n \bar{R}_n \bar{S}_n}{\sqrt{\left(\sum_{i=1}^n R_i^2 - n \bar{R}_n^2 \right) \left(\sum_{i=1}^n S_i^2 - n \bar{S}_n^2 \right)}}. \quad (10.8)$$

Poznámka. Všimněme si, že

$$\bar{R}_n = \bar{S}_n = \frac{n+1}{2}, \quad \sum_{i=1}^n R_i S_i = \frac{1}{2} \sum_{i=1}^n (R_i^2 + S_i^2) - \frac{1}{2} \sum_{i=1}^n (R_i - S_i)^2.$$

Dále za předpokladu, že v datech nejsou shody (tj. všechny hodnoty X_1, \dots, X_n jsou odlišné a také všechny hodnoty Y_1, \dots, Y_n jsou odlišné) platí

$$\sum_{i=1}^n R_i^2 = \sum_{i=1}^n S_i^2 = \sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{6}.$$

Pak lze přepsat $\widehat{\rho}_n^{(S)}$ v jednodušším tvaru:

$$\widehat{\rho}_n^{(S)} = 1 - \frac{6}{n(n^2-1)} \sum_{i=1}^n (R_i - S_i)^2. \quad (10.9)$$

* Angl. *Spearman correlation coefficient*

Jelikož shody nemohou vzniknout, pokud máme náhodný výběr ze spojitého rozdělení, tak někteří autoři předpokládají nejdříve spojitost marginálních rozdělení a definují výběrový Spearmanův korelační koeficient rovnou jako (10.9).

Z definice (10.8) vidíme, že Spearmanův korelační koeficient je vlastně Pearsonův korelační koeficient počítaný z pořadí. Tedy podobně jako u Pearsonova korelačního koeficientu pomocí Cauchy-Schwarzovy nerovnosti dostáváme, že $\widehat{\rho}_n^{(S)} = 1$ právě když existují konstanty $a \in \mathbb{R}$ a $b > 0$ takové, že $R_i = a + bS_i$ pro každé i . Jelikož R_i a S_i jsou pořadí, tak nutně dostáváme, že v tomto případě $R_i = S_i$ (v případě, že v pozorovaných nejsou shody to vlastně okamžitě plyne z přepisu (10.9)). Uvědomme si však, že pořadí X_i a Y_i jsou si rovna, právě když existuje rostoucí funkce h taková, že $Y_i = h(X_i)$ pro každé i . Tedy $\widehat{\rho}_n^{(S)} = 1$, právě když Y_i je rostoucí transformací X_i . Připomeňme, že Pearsonův výběrový korelační koeficient $\widehat{\rho}_n = 1$, právě když X_i je rostoucí lineární transformací Y_i .

Podobně lze odvodit, že $\widehat{\rho}_n^{(S)} = -1$, právě když $R_i = n + 1 - S_i$, tj. když existuje ostře klesající funkce h taková, že $X_i = h(Y_i)$ pro každé i .

Poznámka. Spearmanův korelační koeficient rozhodně není odhadem teoretického korelačního koeficientu ρ daného rovnicí (10.1). Dá se ukázat, že podobně jako výběrový korelační koeficient $\widehat{\rho}_n$ odhaduje teoretický korelační koeficient ρ , tak $\widehat{\rho}_n^{(S)}$ odhaduje

$$\rho^{(S)} = \frac{\text{cov}(F_X(X_i), F_Y(Y_i))}{\sqrt{\text{var}(F_X(X_i)) \text{var}(F_Y(Y_i))}},$$

kde F_X a F_Y jsou distribuční funkce náhodných veličin X_i a Y_i . Tj. na $\rho^{(S)}$ můžeme nahlížet jako na Pearsonův korelační koeficient transformovaných náhodných veličin $F_X(X_i)$ a $F_Y(Y_i)$.

TESTOVÁNÍ NEZÁVISLOSTI

Spearmanův korelační koeficient se zpravidla využívá k testování hypotézy nezávislosti (10.3). Podobně jako pro Pearsonův korelační koeficient se dá ukázat, že pokud X_i a Y_i jsou **nezávislé** a mají nedegenerovaná rozdělení, pak

$$T_n^{(S)} = \sqrt{n-2} \frac{\widehat{\rho}_n^{(S)}}{\sqrt{1 - (\widehat{\rho}_n^{(S)})^2}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1). \quad (10.10)$$

Za povšimnutí stojí, že výše uvedený asymptotický výsledek nevyžaduje (kromě nezávislosti) žádné speciální předpoklady, tj. platí v modelu

$$\mathcal{F} = \{\text{dvourozměrná rozdělení s } \text{var}(X_i) > 0, \text{var}(Y_i) > 0\}.$$

Nyní analogicky jako pro Pearsonův korelační koeficient dostáváme asymptotický kritický obor

$$H_0 \text{ zamítneme} \Leftrightarrow |T_n^{(S)}| \geq t_{n-2}(1 - \alpha/2). \quad (10.11)$$

Poznámka. Pokud v datech nejsou shody, tak se pro malé rozsahy výběru mohou používat přesné kritické hodnoty, které jsou tabelovány.

Podobně jako test založený na Pearsonově korelačním koeficientu v předchozí sekci bude tento test citlivý vůči alternativám (tj. konzistentní pro alternativy), kdy je korelační koeficient mezi $F_X(X_i)$ a $F_Y(Y_i)$ nenulový. Naopak nebude konzistentní proti alternativám, kdy X_i a Y_i sice nejsou nezávislé, ale $\rho^{(S)} = 0$.

Poznámka. V případech, že marginální rozdělení X_i a Y_i jsou spojitá, tak lze ukázat, že za platnosti nulové hypotézy*

$$E \widehat{\rho}_n^{(S)} = 0 \quad \text{a} \quad \text{var}(\widehat{\rho}_n^{(S)}) = \frac{1}{n-1}.$$

Jelikož také za platnosti nezávislosti (v modelu \mathcal{F}) platí

$$\sqrt{n-1} \widehat{\rho}_n^{(S)} \xrightarrow[n \rightarrow \infty]{d} N(0, 1),$$

tak někteří autoři uvádějí kritický obor asymptotického testu nezávislosti ve tvaru

$$H_0 \text{ zamítneme} \Leftrightarrow \sqrt{n-1} |\widehat{\rho}_n^{(S)}| \geq u_{1-\alpha/2}.$$

PORUŠENÍ PŘEDPOKLADŮ

Shody (např. kvůli zaokrouhlování). Kvůli zaokrouhlování bývají v aplikacích v datech často shody. Spearmanův korelační koeficient $\widehat{\rho}_n^{(S)}$ se pak spočte stále dle definice (10.8) s využitím tzv. *průměrných pořadí*, viz. (2.6).[†] V tomto případě stále za nulové hypotézy platí asymptotický výsledek (10.10), tj. lze testovat pomocí kritického oboru (10.11).

Tedy pouze je potřeba si dát pozor na to, že v případě shod **nelze** k výpočtu Spearmanova korelačního koeficientu $\widehat{\rho}_n^{(S)}$ použít vzorec (10.9).

* Viz např. kapitola 8.5 knihy [Hollander et al. \(2013\)](#) † Nutno poznamenat, že literatura v tomto není jednotná. Například v knize [Anděl \(2002\)](#) je Spearmanův korelační koeficient v případě shod definovaný jinak.

Přípravné příklady ke zkoušce.

Vaše řešení „praktických úloh“ by mělo obsahovat matematický model, hypotézu, testovou statistiku a její přesné (či asymptotické) rozdělení za nulové hypotézy. Dále pak kritický obor nebo vzorec pro výpočet p -hodnoty. Mělo by být také řečeno, zda je daný test přesný nebo asymptotický.

1. Máme údaje o tělesné výšce 1000 dospělých mužů a jejich IQ. Navrhněte vhodný test, který by zjistil, zda tělesná výška souvisí s IQ.

A. APPENDIX

A.1. ROZDĚLENÍ χ^2 A t

Definice A.1 (χ^2 -rozdělení) Necht' X_1, \dots, X_k jsou nezávislé náhodné veličiny s rozdělením $N(0, 1)$. Pak rozdělení náhodné veličiny $\sum_{i=1}^k X_i^2$ se nazývá χ^2 -rozdělení o k stupních volnosti. Značíme $Y \sim \chi_k^2$.

Definice A.2 (t -rozdělení) Necht' $X \sim N(0, 1)$ a $Z \sim \chi_k^2$ jsou nezávislé. Pak rozdělení náhodné veličiny $T \stackrel{\text{df}}{=} \frac{X}{\sqrt{Z/k}}$ se nazývá [Studentovo] t -rozdělení s k stupni volnosti. Značíme $T \sim t_k$.

A.2. IDEMPOTENTNÍ MATICE

Definice A.3 Čtvercovou \mathbb{A} (typu $n \times n$) nazveme **idempotentní**, jestliže platí $\mathbb{A}\mathbb{A} = \mathbb{A}$.

Lemma A.1 Vlastními čísly idempotentní matice jsou pouze nuly a jedničky.

Důkaz. Necht' \mathbb{A} je idempotentní matice a λ je vlastní číslo této matice a x příslušný vlastní vektor. Potom z vlastností vlastních čísel a vektorů na jednu stranu platí

$$\mathbb{A}\mathbb{A}x = \mathbb{A}(\mathbb{A}x) = \mathbb{A}\lambda x = \lambda^2 x.$$

Na druhou stranu však z vlastností idempotentní matice platí

$$\mathbb{A}\mathbb{A}x = (\mathbb{A}\mathbb{A})x = \mathbb{A}x = \lambda x.$$

Tedy dostáváme, že $\lambda^2 x = \lambda x$ a tudíž λ může být buď pouze nula nebo jedna. \square

Lemma A.2 Hodnost idempotentní matice se rovná její stopě.

Důkaz. Víme, že pro čtvercovou matici platí, že stopa se rovná součtu vlastních čísel. Dále z lemmatu A.1 víme, že vlastní čísla idempotentní matice jsou pouze nuly a jedničky. Tedy stopa idempotentní matice se rovná násobnosti vlastního čísla jedna. To však odpovídá hodnosti idempotentní matice. \square

A.3. ROZDĚLENÍ KVADRATICKÝCH FOREM

Lemma A.3 Nechť $\mathbf{X} \sim N_d(\mathbf{0}, \Sigma)$. Potom $\sum_{i=1}^d X_i^2$ má stejné rozdělení jako náhodná veličina $\sum_{i=1}^d \lambda_i Y_i^2$, kde Y_1, \dots, Y_d jsou nezávislé náhodné veličiny s rozdělením $N(0, 1)$ a kde $\lambda_1, \dots, \lambda_d$ jsou vlastní čísla matice Σ .

Důkaz. Matice Σ je pozitivně semidefinitní, tudíž dle spektrálního rozkladu existuje ortonormální matice \mathbf{U} taková, že

$$\Sigma = \mathbf{U} \mathbf{D} \mathbf{U}^T$$

kde $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_d)$ je diagonální matice, která má na diagonále vlastní čísla matice Σ .

Z definice mnohorozměrného normálního rozdělení můžeme náhodný vektor \mathbf{X} reprezentovat jako $\mathbf{X} = \mathbf{U} \mathbf{D}^{1/2} \mathbf{Y}$, kde $\mathbf{Y} = (Y_1, \dots, Y_d)^T$ a Y_1, \dots, Y_d jsou nezávislé náhodné veličiny s rozdělením $N(0, 1)$. Tudíž

$$\sum_{i=1}^d X_i^2 = \mathbf{X}^T \mathbf{X} = \mathbf{Y}^T \mathbf{D}^{1/2} \mathbf{U}^T \mathbf{U} \mathbf{D}^{1/2} \mathbf{Y} = \mathbf{Y}^T \mathbf{D} \mathbf{Y} = \sum_{i=1}^d \lambda_i Y_i^2.$$

□

Důležitým speciálním případem je, že $\mathbf{X} \sim N_d(\mathbf{0}, \Sigma)$, kde matice Σ je idempotentní. Potom s využitím lemmat A.1 a A.3 má náhodná veličina $\sum_{i=1}^d X_i^2$ chí-kvadrát rozdělení, kde počet stupňů volnosti odpovídá počtu nenulových vlastních čísel. Ten je však dle lemmatu A.2 roven stopě idempotentní matice. Následující lemma rozšiřuje tuto úvahu na o něco obecnější situaci.

Lemma A.4 Nechť $\mathbf{X} \sim N_n(\mathbf{0}, \Sigma)$ a nechť \mathbb{A} je pozitivně semidefinitní matice typu $n \times n$ taková, že $\mathbb{A}\Sigma$ je nenulová a idempotentní. Pak

$$\mathbf{X}^T \mathbb{A} \mathbf{X} \sim \chi_{\text{tr}(\mathbb{A}\Sigma)}^2.$$

Důkaz. Označme $r = h(\Sigma)$. Z definice mnohorozměrného normálního rozdělení můžeme reprezentovat \mathbf{X} jako $\mathbf{X} = \mathbb{B} \mathbf{Y}$, kde \mathbb{B} je matice typu $n \times r$ taková, že $\Sigma = \mathbb{B} \mathbb{B}^T$ a náhodný vektor \mathbf{Y} má nezávislé složky s rozdělením $N(0, 1)$. Tudíž

$$\mathbf{X}^T \mathbb{A} \mathbf{X} = \mathbf{Y}^T \mathbb{B}^T \mathbb{A} \mathbb{B} \mathbf{Y}.$$

Nyní stačí ukázat, že matice $\mathbb{B}^T \mathbb{A} \mathbb{B}$ je idempotentní. Tvrzení lemmatu pak poplyne z úvah předcházejících toto lemma a dále z toho, že

$$\text{tr}(\mathbb{B}^T \mathbb{A} \mathbb{B}) = \text{tr}(\mathbb{A} \mathbb{B} \mathbb{B}^T) = \text{tr}(\mathbb{A} \Sigma).$$

Dle předpokladu lemmatu je matice $\mathbb{A} \Sigma$ idempotentní, tudíž platí

$$\mathbb{A} \Sigma \mathbb{A} \Sigma = \mathbb{A} \Sigma.$$

Tedy

$$ABB^T ABB^T = ABB^T.$$

Vynásobením výše uvedené rovnice zleva maticí B^T a zprava maticí $(B^T)^{-1}$ pak dává

$$(B^T A B) (B^T A B) = B^T A B,$$

čímž jsme ověřili, že matice $B^T A B$ je vskutku idempotentní. □

A.4. TRANSFORMACE NÁHODNÉ VELIČINY JEJÍ DISTRIBUČNÍ FUNKCÍ

Lemma A.5 Nechť náhodná veličina X má **spojitou** distribuční funkci F . Potom náhodná veličina $F(X)$ má rovnoměrné rozdělení na intervalu $(0, 1)$.

Důkaz. Pro $u \in (0, 1)$ počítejme

$$P[F(X) \leq u] = P[X \leq F^{-1}(u)] = F(F^{-1}(u)) = u,$$

kde v poslední rovnosti jsme využili spojitosti distribuční funkce F . □

Následující lemma je „inverzní“ k předešlému lemmatu a používá se například ke generování náhodných veličin s rozdělením s předepsanou distribuční funkcí F . Za povšimnutí stojí, že oproti lemmatu A.5 není vyžadována spojitost distribuční funkce F .

Lemma A.6 Nechť náhodná veličina U má rovnoměrné rozdělení na intervalu $(0, 1)$ a F je nějaká distribuční funkce. Potom náhodná veličina $F^{-1}(U)$ má rozdělení dané distribuční funkcí F .

Důkaz. Nechť $x \in \mathbb{R}$. Počítejme

$$P[F^{-1}(U) \leq x] = P[U \leq F(x)] = F(x).$$

□

A.5. GAMA FUNKCE A BETA FUNKCE

Gama funkce je definována jako

$$\Gamma(z) = \int_0^{\infty} x^{z-1} e^{-x} dx,$$

přičemž pro naše potřeby postačí uvažovat $z > 0$. Z vlastností gama funkce se v příkladech často využívá, že pro $n \in \mathbb{N}$ platí $\Gamma(n) = (n-1)!$.

Beta funkce je pro $a, b > 0$ definována jako

$$B(a, b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx.$$

Platí

$$B(a, b) = \frac{\Gamma(a) \Gamma(b)}{\Gamma(a+b)}. \quad (\text{A.1})$$

LITERATURA

- Anděl, J. (1985). *Matematická statistika*. SNTL, Praha.
- Anděl, J. (1998). *Statistické metody*. Praha: Matfyzpress.
- Anděl, J. (2002). *Základy matematické statistiky*. Praha: Matfyzpress.
- Chung, E. and J. P. Romano (2016). Asymptotically valid and exact permutation tests based on two-sample U-statistics. *Journal of Statistical Planning and Inference* 168, 97–105.
- Cochran, W. G. (1954). Some methods for strengthening the common χ^2 -tests. *Biometrics* 10(4), 417–451.
- Dupač, V. and M. Hušková (1999). *Pravděpodobnost a matematická statistika*. Praha: Karolinum.
- Hájek, J. and Z. Šidák (1967). *Theory of rank tests*. Academia, Praha.
- Hollander, M., D. A. Wolfe, and E. Chicken (2013). *Nonparametric statistical methods*. John Wiley & Sons, New York.
- Kika, V. (2015). Analýza rozptylu při nesplnění předpokladu normality. Bakalářská práce. Univerzita Karlova, Matematicko-fyzikální fakulta.
- Kulich, M. (2018). Základy teorie pravděpodobnosti pro předmět Matematická statistika 1. https://www.karlin.mff.cuni.cz/~kulich/vyuka/ms1/doc/pravdepodobnost_ms1.pdf.
- Lachout, P. (2004). *Teorie pravděpodobnosti*. Karolinum. Skripta.
- Serfling, R. J. (1980). *Approximation Theorems of Mathematical Statistics*. Wiley, New York.
- Welch, B. L. (1938). The significance of the difference between two means when the population variances are unequal. *Biometrika* 29(3/4), 350–362.
- Welch, B. L. (1951). On the comparison of several mean values: An alternative approach. *Biometrika* 38, 330–336.