

Pravděpodobnost a matematická statistika

Mirko Navara
Centrum strojového vnímání
katedra kybernetiky FEL ČVUT
Karlovo náměstí, budova G, místnost 104a
<http://cmp.felk.cvut.cz/~navara/MVT>
<http://cmp.felk.cvut.cz/~navara/psi>

7. února 2013

Obsah

1	O čem to je?	4
1.1	Teorie pravděpodobnosti	4
1.2	Statistika	4
2	Základní pojmy teorie pravděpodobnosti	4
2.1	Laplaceova (klasická) definice pravděpodobnosti	4
2.1.1	Základní pojmy	5
2.1.2	Pravděpodobnost	5
2.1.3	Náhodná veličina	5
2.2	Vlastnosti pravděpodobnosti	5
2.2.1	Úplný systém jevů	5
2.3	Problémy Laplaceovy definice pravděpodobnosti	6
2.3.1	Rozšíření Laplaceova modelu pravděpodobnosti	6
2.4	Kombinatorické pojmy a vzorce	6
2.5	Kolmogorovova definice pravděpodobnosti	8
2.5.1	Borelova σ-algebra	8
2.5.2	Pravděpodobnost (=pravděpodobnostní míra)	8
3	Nezávislost a podmíněná pravděpodobnost	9
3.1	Nezávislé jevy	9
3.2	Podmíněná pravděpodobnost	10
3.2.1	Podmíněná nezávislost	11
4	Náhodné veličiny a vektory	11
4.1	Náhodná veličina	12
4.2	n -rozměrný náhodný vektor (n -rozměrná náhodná veličina)	13
4.3	Nezávislost náhodných veličin	14
4.4	Obecnější náhodné veličiny	14

4.5	Směs náhodných veličin	15
4.6	Druhy náhodných veličin	16
4.7	Popis spojité náhodné veličiny	16
4.8	Popis smíšené náhodné veličiny	17
4.9	Kvantilová funkce náhodné veličiny	17
4.10	Jak reprezentovat náhodnou veličinu v počítači	18
4.11	Operace s náhodnými veličinami	19
4.12	Jak realizovat náhodnou veličinu na počítači	20
4.13	Střední hodnota	20
4.13.1	Vlastnosti střední hodnoty	21
4.14	Rozptyl (disperze)	22
4.15	Směrodatná odchylka	22
4.16	Obecné a centrální momenty	23
4.17	Normovaná náhodná veličina	23
4.18	Základní typy diskrétních rozdělení	23
4.18.1	Diracovo	23
4.18.2	Rovnoměrné	24
4.18.3	Alternativní (Bernoulliovo)	24
4.18.4	Binomické $Bi(m, q)$	24
4.18.5	Poissonovo $Po(\lambda)$	24
4.18.6	Geometrické	25
4.18.7	Hypergeometrické	25
4.19	Základní typy spojitých rozdělení	26
4.19.1	Rovnoměrné $R(a, b)$	26
4.19.2	Normální (Gaussovo) $N(\mu, \sigma^2)$	26
4.19.3	Logaritmicnormální $LN(\mu, \sigma^2)$	26
4.19.4	Exponenciální $Ex(\tau)$	27
4.20	Náhodné vektory 2	27
4.20.1	Diskrétní náhodný vektor	27
4.20.2	Spojité náhodný vektor	27
4.21	Číselné charakteristiky náhodného vektoru	28
4.21.1	Vícerozměrné normální rozdělení $N(\mu, \Sigma)$	30
4.22	Lineární prostor náhodných veličin	30
4.22.1	Lineární podprostor \mathcal{N} náhodných veličin s nulovými středními hodnotami	31
4.22.2	Lineární regrese	31
4.23	Reprezentace náhodných vektorů v počítači	32
4.24	Čebyševova nerovnost	32

5 Základní pojmy statistiky 33

5.1	K čemu potřebujeme statistiku	33
5.2	Pojem náhodného výběru, odhady	33
5.3	Výběrový průměr	34
5.4	Výběrový rozptyl	35
5.4.1	Rozdělení χ^2 s 1 stupněm volnosti	36
5.4.2	Rozdělení χ^2 s η stupni volnosti	37
5.4.3	Výběrový rozptyl	37
5.4.4	Alternativní odhad rozptylu	38

5.5	Výběrová směrodatná odchylka	39
5.6	Výběrový k -tý obecný moment	39
5.7	Histogram a empirické rozdělení	40
5.7.1	Vlastnosti empirického rozdělení	40
5.8	Výběrový medián	40
5.9	Intervalové odhady	41
5.10	Intervalové odhady parametrů normálního rozdělení $N(\mu, \sigma^2)$	41
5.10.1	Odhad střední hodnoty při známém rozptylu σ^2	41
5.10.2	Odhad střední hodnoty při neznámém rozptylu	42
5.10.3	Studentovo t-rozdělení (autor: Gossett)	42
5.10.4	Odhad střední hodnoty při neznámém rozptylu II	43
5.10.5	Odhad rozptylu	43
5.10.6	Intervalové odhady spojitých rozdělení, která nejsou normální	44
5.11	Obecné odhady parametrů	44
5.11.1	Metoda momentů	44
5.11.2	Metoda maximální věrohodnosti (likelihood)	45

6 Testování hypotéz 49

6.1	Základní pojmy a principy testování hypotéz	49
6.2	Testy střední hodnoty normálního rozdělení	51
6.2.1	Při známém rozptylu σ^2	51
6.2.2	Při neznámém rozptylu	51
6.3	Testy rozptylu normálního rozdělení	52
6.4	Porovnání dvou normálních rozdělení	52
6.4.1	Testy rozptylu dvou normálních rozdělení [Fisher]	52
6.4.2	Testy středních hodnot dvou normálních rozdělení se známým rozptylem σ^2	53
6.4.3	Testy středních hodnot dvou normálních rozdělení se (stejným) neznámým rozptylem	53
6.5	Testy středních hodnot dvou normálních rozdělení - párový pokus	54
6.5.1	Pro známý rozptyl σ^2	55
6.5.2	Pro neznámý rozptyl	55
6.6	χ^2 -test dobré shody	56
6.6.1	Modifikace	57
6.6.2	χ^2 -test dobré shody dvou rozdělení	57
6.6.3	χ^2 -test nezávislosti dvou rozdělení	58
6.7	Korelace, její odhad a testování	58
6.7.1	Test nekorelovanosti dvou normálních rozdělení	59
6.8	Neparametrické testy	59
6.8.1	Znaménkový test	59
6.8.2	Wilcoxonův test (jednovýběrový)	60

7 Co zde nebylo 60

7.1	Více o zobrazení náhodné veličiny funkcí a o součtu náhodných veličin	60
7.2	Diskretizace	60
7.3	Směs pravděpodobností	60
7.4	Charakteristická funkce náhodné veličiny	60
7.5	Důkaz centrální limitní věty	60

1 O čem to je?

1A. Jak vysoká by měla být pojistka auta proti krádeži (bez marže), je-li jeho cena 1 000 000 Kč a riziko ukradení během pojistného období 0.001?

$$1\,000\,000 \cdot 0.001 = 1\,000 \text{ Kč}$$

1B. Jak vysoká by měla být pojistka auta pro případ havárie, při níž může být škoda různě velká?

⇒ **TEORIE PRAVDĚPODOBNOTI**

2. Jak odhadnout pravděpodobnost krádeže auta nebo střední škodu při havárii a jak přesný bude odhad?

⇒ **STATISTIKA**

3. Jak označovat auta a jejich díly, abychom je jednoznačně určili?

⇒ **TEORIE INFORMACE A KÓDOVÁNÍ**

Určitě ne jako v rozvrhu na FEL: Cvičení AD0B01PSI bude v učebně KN:E-24.

1.1 Teorie pravděpodobnosti

je nástroj pro účelné rozhodování v systémech, kde **budoucí** pravdivost jevů závisí na okolnostech, které zcela neznáme.

Poskytuje model takových systémů a kvantifikaci výsledků.

Pravděpodobnostní **popis** ⇒ **chování** systému

1.2 Statistika

je nástroj pro hledání a ověřování pravděpodobnostního popisu reálných systémů na základě jejich pozorování.

Chování systému ⇒ pravděpodobnostní **popis**

Poskytuje daleko víc: nástroj pro zkoumání světa, pro hledání a ověřování závislostí, které nejsou zjevné.

2 Základní pojmy teorie pravděpodobnosti

2.1 Laplaceova (klasická) definice pravděpodobnosti

Předpoklad: Náhodný pokus s $n \in \mathbb{N}$ různými, vzájemně se vylučujícími výsledky, které jsou **stejně možné**.

Pravděpodobnost jevu, který nastává právě při k z těchto výsledků, je k/n .

1. problém: „stejně možné“=„stejně pravděpodobné,“ ale co to znamená? (definice kruhem!)

Elementární jevy jsou všechny „stejně možné“ výsledky.

Množina všech elementárních jevů: Ω

Jev: $A \subseteq \Omega$

Úmluva. Jevy budeme ztotožňovat s příslušnými množinami elementárních jevů a používat pro ně množinové operace (místo výrokových).

2.1.1 Základní pojmy

Jev jistý: $\Omega, \mathbf{1}$

Jev nemožný: $\emptyset, \mathbf{0}$

Konjunkce jevů („and“): $A \cap B$

Disjunkce jevů („or“): $A \cup B$

Jev opačný k A : $\bar{A} = \Omega \setminus A$

$A \Rightarrow B: A \subseteq B$

Jevy neslučitelné (=vzájemně se vylučující): $A_1, \dots, A_n: \bigcap_{i \leq n} A_i = \emptyset$

Jevy po dvou neslučitelné: $A_1, \dots, A_n: \forall i, j \in \{1, \dots, n\}, i \neq j: A_i \cap A_j = \emptyset$

Jevové pole: všechny jevy pozorovatelné v náhodném pokusu, zde $\exp \Omega$ (=množina všech podmnožin množiny Ω)

2.1.2 Pravděpodobnost

jevu A :

$$P(A) = \frac{|A|}{|\Omega|},$$

kde $|\cdot|$ značí počet prvků množiny

2.1.3 Náhodná veličina

je libovolná funkce $X: \Omega \rightarrow \mathbb{R}$

Střední hodnota:

$$EX = \frac{1}{n} \sum_{\omega \in \Omega} X(\omega),$$

kde $n = |\Omega|$.

Příklad: Elementární jevy jsou možné výsledky hry, náhodná veličina je výše výhry. Střední hodnota je spravedlivá cena za účast ve hře.

2.2 Vlastnosti pravděpodobnosti

$$P(A) \in \langle 0, 1 \rangle$$

$$P(\mathbf{0}) = 0, \quad P(\mathbf{1}) = 1$$

$$P(\bar{A}) = 1 - P(A)$$

$$A \subseteq B \Rightarrow P(A) \leq P(B)$$

$$A \subseteq B \Rightarrow P(B \setminus A) = P(B) - P(A)$$

$$A \cap B = \emptyset \Rightarrow P(A \cup B) = P(A) + P(B) \quad (\text{aditivita})$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

2.2.1 Úplný systém jevů

tvorí jevy $B_i, i \in I$, jestliže jsou po dvou neslučitelné a $\bigcup_{i \in I} B_i = \mathbf{1}$.

Speciální případ pro 2 jevy: $\{C, \bar{C}\}$

Je-li $\{B_1, \dots, B_n\}$ **úplný systém jevů**, pak

$$\sum_{i=1}^n P(B_i) = 1$$

a pro libovolný jev A

$$P(A) = \sum_{i=1}^n P(A \cap B_i).$$

Speciálně:

$$P(A) = P(A \cap C) + P(A \cap \overline{C}).$$

2.3 Problémy Laplaceovy definice pravděpodobnosti

2. problém: Nedovoluje nekonečné množiny jevů, geometrickou pravděpodobnost...

Nelze mít nekonečně mnoho stejně pravděpodobných výsledků.

Příklad: Podíl plochy pevniny k povrchu Země je pravděpodobnost, že náhodně vybraný bod na Zemi leží na pevnině (je-li výběr bodů prováděn „rovnoměrně“).

Příklad (Buffonova úloha): Na linkovaný papír hodíme jehlu, jejíž délka je rovna vzdálenosti mezi linkami. Jaká je pravděpodobnost, že jehla protne nějakou linku?

3. problém: Nedovoluje iracionální hodnoty pravděpodobnosti.

2.3.1 Rozšíření Laplaceova modelu pravděpodobnosti

Příklad: Místo hrací kostky házíme krabičkou od zápalek, jejíž strany jsou nestejně dlouhé. Jaká je pravděpodobnost možných výsledků?

Připustíme, že **elementární jevy nemusí být stejně pravděpodobné**.

Ztrácíme návod, jak pravděpodobnost stanovit. Je to funkce, která jevům přiřazuje čísla z intervalu $\langle 0, 1 \rangle$ a splňuje jisté podmínky. Nemáme návod, jak z nich vybrat tu pravou.

Tato nevýhoda je neodstranitelná a je důvodem pro vznik statistiky, která k danému opakovatelnému pokusu hledá pravděpodobnostní model.

2.4 Kombinatorické pojmy a vzorce

(Dle [Zvára, Štěpán].)

V urně je n rozlišitelných objektů, postupně vytáhneme k .

výběr	s vracením	bez vracení
uspořádaný	variace s opakováním n^k	variace bez opakování $\frac{n!}{(n-k)!}$
neuspořádaný	kombinace s opakováním $\binom{n+k-1}{k}$	kombinace bez opakování $\frac{n!}{k!(n-k)!} = \binom{n}{k}$

Z této tabulky pouze **kombinace s opakováním** nejsou všechny stejně pravděpodobné (odpovídají různému počtu variací s opakováním) a nedovolují proto použití Laplaceova modelu pravděpodobnosti.

Permutace (pořadí) bez opakování: Tvoříme posloupnost z n hodnot, přičemž každá se vyskytne právě jednou. Počet permutací je $n!$ (je to speciální případ variací bez opakování pro $n = k$).

Permutace s opakováním: Tvoříme posloupnost délky k z n hodnot, přičemž j -tá hodnota se opakuje k_j -krát, $\sum_{j=1}^n k_j = k$. Počet **různých** posloupností je

$$\frac{k!}{k_1! \cdot \dots \cdot k_n!}.$$

Speciálně pro $n = 2$ dostáváme

$$\frac{k!}{k_1! \cdot k_2!} = \frac{k!}{k_1! \cdot (k - k_1)!} = \binom{k}{k_1},$$

což je počet kombinací bez opakování (ovšem k_1 -prvkových z k prvků).

n	4	10	100	1 000	10 000
počet 4-prvkových variací z n prvků bez opakování, $\frac{n!}{(n-4)!}$	24	5 040	94 109 400	$0.994 \cdot 10^{12}$	$0.9994 \cdot 10^{16}$
počet 4-prvkových variací z n prvků s opakováním, n^4	256	10 000	10^8	10^{12}	10^{16}
počet 4-prvkových kombinací z n prvků bez opakování, $\binom{n}{4}$	1	210	3 921 225	41 417 124 750	$4.164 \cdot 10^{14}$
počet 4-prvkových kombinací z n prvků s opakováním, $\binom{n+3}{4}$	35	715	4 421 275	41 917 125 250	$4.169 \cdot 10^{14}$

Věta 1. Pro dané $k \in \mathbb{N}$ a pro $n \rightarrow \infty$ se poměr počtů variací (resp. kombinací) bez opakování a s opakováním blíží jedné, tj.

$$\lim_{n \rightarrow \infty} \frac{n!}{(n-k)! n^k} = 1, \quad \lim_{n \rightarrow \infty} \frac{\binom{n}{k}}{\binom{n+k-1}{k}} = 1.$$

Důkaz.

$$\begin{aligned} \frac{n!}{(n-k)! n^k} &= \frac{n(n-1) \cdots (n-(k-1))}{n^k} = \\ &= 1 \left(1 - \frac{1}{n}\right) \cdots \left(1 - \frac{k-1}{n}\right) \rightarrow 1, \\ \frac{\binom{n}{k}}{\binom{n+k-1}{k}} &= \frac{n(n-1) \cdots (n-(k-1))}{(n+(k-1)) \cdots (n+1)n} = \\ &= \frac{1 \left(1 - \frac{1}{n}\right) \cdots \left(1 - \frac{k-1}{n}\right)}{\left(1 + \frac{k-1}{n}\right) \cdots \left(1 + \frac{1}{n}\right) 1} \rightarrow 1 \end{aligned}$$

(počet činitelů k je konstantní). □

Důsledek 1. Pro $n \gg k$ je počet variací (resp. kombinací) s opakováním přibližně

$$\frac{n!}{(n-k)!} \doteq n^k, \quad \binom{n}{k} \doteq \frac{n^k}{k!}.$$

Jednodušší bývá **neuspořádaný výběr bez vracení** nebo **uspořádaný výběr s vracením**.

2.5 Kolmogorovova definice pravděpodobnosti

Elementárních jevů (=prvků množiny Ω) může být **nekonečně mnoho, nemusí být stejně pravděpodobné**.

Jevy jsou podmnožiny množiny Ω , ale **ne nutně všechny**; tvoří podmnožinu $\mathcal{A} \subseteq \exp \Omega$, která splňuje následující podmínky:

$$(A1) \quad \emptyset \in \mathcal{A}.$$

$$(A2) \quad A \in \mathcal{A} \Rightarrow \bar{A} \in \mathcal{A}.$$

$$(A3) \quad (\forall n \in \mathbb{N} : A_n \in \mathcal{A}) \Rightarrow \bigcup_{n \in \mathbb{N}} A_n \in \mathcal{A}.$$

Systém \mathcal{A} podmnožin nějaké množiny Ω , který splňuje podmínky (A1-3), se nazývá **σ -algebra**.

Důsledky: $\Omega = \bar{\emptyset} \in \mathcal{A}$,

$$(\forall n \in \mathbb{N} : A_n \in \mathcal{A}) \Rightarrow \bigcap_{n \in \mathbb{N}} A_n = \overline{\bigcup_{n \in \mathbb{N}} \bar{A}_n} \in \mathcal{A}.$$

Přirozený nápad $\mathcal{A} = \exp \Omega$ vede k nežádoucím paradoxům.

(A3) je uzavřenost na **spočetná** sjednocení.

Uzavřenost na **jakákoli** sjednocení se ukazuje jako příliš silný požadavek.

Uzavřenost na **konečná** sjednocení se ukazuje jako příliš slabý požadavek; nedovoluje např. vyjádřit kruh jako sjednocení obdélníků.

\mathcal{A} nemusí ani obsahovat všechny jednobodové množiny, v tom případě **elementární jevy nemusí být jevy!**

2.5.1 Borelova σ -algebra

je nejmenší σ -algebra podmnožin \mathbb{R} , která obsahuje všechny intervaly.

Obsahuje všechny intervaly otevřené, uzavřené i polouzavřené, i jejich spočetná sjednocení, a některé další množiny, ale je menší než $\exp \mathbb{R}$. Její prvky nazýváme **borelovské množiny**.

2.5.2 Pravděpodobnost (=pravděpodobnostní míra)

je funkce $P: \mathcal{A} \rightarrow \langle 0, 1 \rangle$, splňující podmínky

$$(P1) \quad P(\mathbf{1}) = 1,$$

$$(P2) \quad P\left(\bigcup_{n \in \mathbb{N}} A_n\right) = \sum_{n \in \mathbb{N}} P(A_n), \text{ pokud jsou množiny (=jevy) } A_n, n \in \mathbb{N}, \text{ po dvou neslučitelné.}$$

(**spočetná aditivita**)

Pravděpodobnostní prostor je trojice (Ω, \mathcal{A}, P) , kde Ω je neprázdná množina, \mathcal{A} je σ -algebra podmnožin množiny Ω a $P: \mathcal{A} \rightarrow \langle 0, 1 \rangle$ je pravděpodobnost.

Dříve uvedené vlastnosti pravděpodobnosti jsou důsledkem (P1), (P2).

(**Konečná**) **aditivita** by byla příliš slabá, nedovoluje např. přechod od obsahu obdélníka k obsahu kruhu.

Příklad („nekonečná ruleta“): Výsledkem může být libovolné přirozené číslo, každé má pravděpodobnost 0.

Úplná aditivita (pro jakékoli soubory po dvou neslučitelných jevů) by byla příliš silným požadavkem. Pak bychom nepřipouštěli ani rovnoměrné rozdělení na intervalu nebo na ploše. Pravděpodobnost zachovává limity monotónních posloupností jevů (množin):
Nechť $(A_n)_{n \in \mathbb{N}}$ je posloupnost jevů.

$$A_1 \subseteq A_2 \subseteq \dots \Rightarrow P\left(\bigcup_{n \in \mathbb{N}} A_n\right) = \lim_{n \rightarrow \infty} P(A_n),$$

$$A_1 \supseteq A_2 \supseteq \dots \Rightarrow P\left(\bigcap_{n \in \mathbb{N}} A_n\right) = \lim_{n \rightarrow \infty} P(A_n).$$

Laplaceův model	Kolmogorovův model
konečně mnoho jevů	i nekonečně mnoho jevů
p-sti jen racionální	p-sti iracionální
$P(A) = 0 \Rightarrow A = \mathbf{0}$	možné jevy s nulovou p-stí
p-sti určeny strukturou jevů	p-sti neurčeny strukturou jevů

3 Nezávislost a podmíněná pravděpodobnost

3.1 Nezávislé jevy

Motivace: Dva jevy spolu „nesouvisí“

Definice: $P(A \cap B) = P(A) \cdot P(B)$.

To je ovšem jen náhražka, která říká mnohem méně, než jsme chtěli!
(Podobně jako $P(A \cap B) = 0$ neznamená, že jevy A, B jsou neslučitelné.)
Pro nezávislé jevy A, B

$$P(A \cup B) = P(A) + P(B) - P(A) \cdot P(B).$$

Důkaz:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) = P(A) + P(B) - P(A) \cdot P(B).$$

Jsou-li jevy A, B nezávislé, pak jsou nezávislé také jevy A, \bar{B} (a též dvojice jevů \bar{A}, B a \bar{A}, \bar{B}).

Důkaz:

$$\begin{aligned} P(A \cap \bar{B}) &= P(A) - P(A \cap B) = P(A) - P(A) \cdot P(B) = \\ &= P(A) \cdot (1 - P(B)) = P(A) \cdot P(\bar{B}). \end{aligned}$$

Jevy A_1, \dots, A_n se nazývají **po dvou nezávislé**, jestliže každé dva z nich jsou nezávislé.
To je málo.

Množina jevů \mathcal{M} se nazývá **nezávislá**, jestliže

$$P\left(\bigcap_{A \in \mathcal{K}} A\right) = \prod_{A \in \mathcal{K}} P(A)$$

pro všechny **konečné** podmnožiny $\mathcal{K} \subseteq \mathcal{M}$.

3.2 Podmíněná pravděpodobnost

Příklad: Fotbalová družstva mohla mít před zápasem rovné šance na vítězství. Je-li však stav zápasu 5 minut před koncem 3 : 0, pravděpodobnosti výhry jsou jiné.

Máme pravděpodobnostní popis systému. Dostaneme-li dodatečnou informaci, že nastal jev B , aktualizujeme naši znalost o pravděpodobnosti jevu A na

$$P(A|B) = \frac{P(A \cap B)}{P(B)},$$

což je **podmíněná pravděpodobnost** jevu A za podmínky B . Je definována pouze pro $P(B) \neq 0$. (To předpokládáme i nadále.)

V novém modelu je $P(\bar{B}|B) = 0$, což odráží naši znalost, že jev \bar{B} nenastal.

Podmíněná pravděpodobnost je chápána též jako funkce

$$P(\cdot|B): \mathcal{A} \rightarrow \langle 0, 1 \rangle, \quad A \mapsto \frac{P(A \cap B)}{P(B)}$$

a je to pravděpodobnost v původním smyslu.

Vlastnosti podmíněné pravděpodobnosti:

- $P(\mathbf{1}|B) = 1$, $P(\mathbf{0}|B) = 0$.
- Jsou-li jevy A_1, A_2, \dots jsou po dvou neslučitelné, pak

$$P\left(\bigcup_{n \in \mathbb{N}} A_n \mid B\right) = \sum_{n \in \mathbb{N}} P(A_n | B).$$

- Je-li $P(A|B)$ definována, jsou jevy A, B **nezávislé**, právě když $P(A|B) = P(A)$.
- $B \subseteq A \Rightarrow P(A|B) = 1$, $P(A \cap B) = 0 \Rightarrow P(A|B) = 0$.

Věta o úplné pravděpodobnosti: Necht' $B_i, i \in I$, je (spočetný) úplný systém jevů a $\forall i \in I : P(B_i) \neq 0$. Pak pro každý jev A platí

$$P(A) = \sum_{i \in I} P(B_i) P(A|B_i).$$

Důkaz:

$$\begin{aligned} P(A) &= P\left(\left(\bigcup_{j \in I} B_j\right) \cap A\right) = P\left(\bigcup_{j \in I} (B_j \cap A)\right) = \\ &= \sum_{i \in I} P(B_i \cap A) = \sum_{i \in I} P(B_i) P(A|B_i). \end{aligned}$$

Příklad: Test nemoci je u 1% zdravých falešně pozitivní a u 10% nemocných falešně negativní. Nemocných je v populaci 0.001. Jaká je pravděpodobnost, že pacient s pozitivním testem je nemocný?

Bayesova věta: Necht' $B_i, i \in I$, je (spočetný) úplný systém jevů a $\forall i \in I : P(B_i) \neq 0$. Pak pro každý jev A splňující $P(A) \neq 0$ platí

$$P(B_i|A) = \frac{P(B_i) P(A|B_i)}{\sum_{j \in I} P(B_j) P(A|B_j)}.$$

Důkaz (s využitím věty o úplné pravděpodobnosti):

$$P(B_i|A) = \frac{P(B_i \cap A)}{P(A)} = \frac{P(B_i) P(A|B_i)}{\sum_{j \in I} P(B_j) P(A|B_j)}.$$

Význam: Pravděpodobnosti $P(A|B_i)$ odhadneme z pokusů nebo z modelu, pomocí nich určíme pravděpodobnosti $P(B_i|A)$, které slouží k „optimálnímu“ odhadu, který z jevů B_i nastal.

Problém: Ke stanovení **aposteriorní pravděpodobnosti** $P(B_i|A)$ potřebujeme znát i **apriorní pravděpodobnost** $P(B_i)$.

Příklad: Na vstupu informačního kanálu mohou být znaky $1, \dots, m$, výskyt znaku j označujeme jako jev B_j . Na výstupu mohou být znaky $1, \dots, k$, výskyt znaku i označujeme jako jev A_i . (Obvykle $k = m$, ale není to nutné.) Obvykle lze odhadnout podmíněné pravděpodobnosti $P(A_i|B_j)$, že znak j bude přijat jako i . Pokud známe apriorní pravděpodobnosti (vyslání znaku j) $P(B_j)$, můžeme pravděpodobnosti příjmu znaků vypočítat maticovým násobením:

$$\begin{aligned} & [P(A_1) \quad P(A_2) \quad \cdots \quad P(A_k)] = \\ & = [P(B_1) \quad P(B_2) \quad \cdots \quad P(B_m)] \cdot \begin{bmatrix} P(A_1|B_1) & P(A_2|B_1) & \cdots & P(A_k|B_1) \\ P(A_1|B_2) & P(A_2|B_2) & \cdots & P(A_k|B_2) \\ \vdots & \vdots & \ddots & \vdots \\ P(A_1|B_m) & P(A_2|B_m) & \cdots & P(A_k|B_m) \end{bmatrix}. \end{aligned}$$

Všechny matice v tomto vzorci mají jednotkové součty řádků (takové matice nazýváme **stochastické**). Podmíněné rozdělení pravděpodobnosti, pokud byl přijat znak i , je

$$P(B_j|A_i) = \frac{P(A_i|B_j) P(B_j)}{P(A_i)}.$$

Rozdělení pravděpodobností vyslaných znaků je

$$\begin{aligned} & [P(B_1) \quad P(B_2) \quad \cdots \quad P(B_m)] = \\ & = [P(A_1) \quad P(A_2) \quad \cdots \quad P(A_k)] \cdot \begin{bmatrix} P(A_1|B_1) & P(A_2|B_1) & \cdots & P(A_k|B_1) \\ P(A_1|B_2) & P(A_2|B_2) & \cdots & P(A_k|B_2) \\ \vdots & \vdots & \ddots & \vdots \\ P(A_1|B_m) & P(A_2|B_m) & \cdots & P(A_k|B_m) \end{bmatrix}^{-1}, \end{aligned}$$

pokud $k = m$ a příslušná inverzní matice existuje.

3.2.1 Podmíněná nezávislost

Náhodné jevy A, B jsou **podmíněně nezávislé** za podmínky C , jestliže

$$P(A \cap B|C) = P(A|C) P(B|C).$$

Podobně definujeme podmíněnou nezávislost více jevů.

4 Náhodné veličiny a vektory

Příklad: Auto v ceně 10 000 EUR bude do roka ukradeno s pravděpodobností 1 : 1 000. Adekvátní cena ročního pojistného (bez zisku pojišťovny) je 10 000/1 000 = 10 EUR.

Někdy tento jednoduchý postup selhává:

Příklad: Pro stanovení havarijního pojištění potřebujeme znát nejen pravděpodobnost havárie (resp. počtu havárií za pojistné období), ale i „průměrnou“ škodu při jedné havárii, lépe pravděpodobnostní rozdělení výše škody.

⇒ Musíme studovat i náhodné pokusy, jejichž výsledky nejsou jen dva (jev nastal/nenastal), ale více hodnot, vyjádřených reálnými čísly.

4.1 Náhodná veličina

na pravděpodobnostním prostoru (Ω, \mathcal{A}, P) je **měřitelná** funkce $X: \Omega \rightarrow \mathbb{R}$, tj. taková, že pro každý interval I platí

$$X^{-1}(I) = \{\omega \in \Omega \mid X(\omega) \in I\} \in \mathcal{A}$$

Je popsána pravděpodobnostmi

$$P_X(I) = P[X \in I] = P(\{\omega \in \Omega \mid X(\omega) \in I\}),$$

definovanými pro libovolný interval I (a tedy i pro libovolné sjednocení spočetně mnoha intervalů a pro libovolnou borelovskou množinu).

P_X je **pravděpodobnostní míra** na Borelově σ -algebře určující **rozdělení náhodné veličiny** X .

K tomu, aby stačila znalost P_X na intervalech, se potřebujeme omezit na tzv. *perfektní míry*; s jinými se v praxi neseťkáme.

Pravděpodobnostní míra P_X splňuje podmínky:

$$P_X(\mathbb{R}) = 1,$$

$$P_X\left(\bigcup_{n \in \mathbb{N}} I_n\right) = \sum_{n \in \mathbb{N}} P_X(I_n), \text{ pokud jsou množiny } I_n, n \in \mathbb{N}, \text{ navzájem disjunktní.}$$

Z toho vyplývá:

$$P_X(\emptyset) = 0, \quad P_X(\mathbb{R} \setminus I) = 1 - P_X(I),$$

$$\text{jestliže } I \subseteq J, \text{ pak } P_X(I) \leq P_X(J) \text{ a } P_X(J \setminus I) = P_X(J) - P_X(I).$$

Úspornější reprezentace: omezíme se na intervaly tvaru $I = (-\infty, t]$, $t \in \mathbb{R}$,

$$P[X \in (-\infty, t]] = P[X \leq t] = P_X((-\infty, t]) = F_X(t).$$

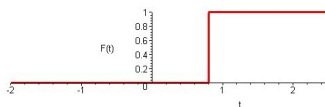
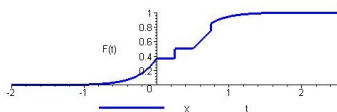
$F_X: \mathbb{R} \rightarrow \langle 0, 1 \rangle$ je **distribuční funkce** náhodné veličiny X . Ta stačí, neboť

$$\begin{aligned} (a, b) &= (-\infty, b] \setminus (-\infty, a], & P_X((a, b)) &= P[a < X \leq b] = F_X(b) - F_X(a), \\ (a, \infty) &= \mathbb{R} \setminus (-\infty, a], & P_X((a, \infty)) &= 1 - F_X(a), \\ (-\infty, a) &= \bigcup_{b: b < a} (-\infty, b], & P_X((-\infty, a)) &= P[X < a] = \lim_{b \rightarrow a^-} F_X(b) = F_X(a-), \\ \{a\} &= (-\infty, a] \setminus (-\infty, a), & P_X(\{a\}) &= P[X = a] = F_X(a) - F_X(a-), \\ \dots & & \dots & \end{aligned}$$

Vlastnosti distribuční funkce:

- neklesající,
- zprava spojitá,
- $\lim_{t \rightarrow -\infty} F_X(t) = 0, \quad \lim_{t \rightarrow \infty} F_X(t) = 1.$

Věta: Tyto podmínky jsou nejen **nutné**, ale i **postačující**.



Příklad: Reálnému číslu r odpovídá náhodná veličina (značená též r) s **Diracovým** rozdělením v r :

$$P_r(I) = \begin{cases} 0 & \text{pro } r \notin I, \\ 1 & \text{pro } r \in I, \end{cases} \quad F_r(t) = \begin{cases} 0 & \text{pro } t < r, \\ 1 & \text{pro } t \geq r. \end{cases}$$

(F_r je posunutá Heavisideova funkce.)

Tvrzení: $X \leq Y \Rightarrow F_X \geq F_Y$.

4.2 n -rozměrný náhodný vektor (n -rozměrná náhodná veličina)

na pravděpodobnostním prostoru (Ω, \mathcal{A}, P) je **měřitelná** funkce $\mathbf{X}: \Omega \rightarrow \mathbb{R}^n$, tj. taková, že pro každý n -rozměrný interval I platí

$$\mathbf{X}^{-1}(I) = \{\omega \in \Omega \mid \mathbf{X}(\omega) \in I\} \in \mathcal{A}.$$

Lze psát

$$\mathbf{X}(\omega) = (X_1(\omega), \dots, X_n(\omega)),$$

kde zobrazení $X_k: \Omega \rightarrow \mathbb{R}$, $k = 1, \dots, n$, jsou náhodné veličiny.

Náhodný vektor lze považovat za vektor náhodných veličin $\mathbf{X} = (X_1, \dots, X_n)$.

Je popsán pravděpodobnostmi

$$\begin{aligned} P_{\mathbf{X}}(I_1 \times \dots \times I_n) &= P[X_1 \in I_1, \dots, X_n \in I_n] = \\ &= P(\{\omega \in \Omega \mid X_1(\omega) \in I_1, \dots, X_n(\omega) \in I_n\}), \end{aligned}$$

kde I_1, \dots, I_n jsou intervaly v \mathbb{R} .

Z těch vyplývají pravděpodobnosti

$$P_{\mathbf{X}}(I) = P[\mathbf{X} \in I] = P(\{\omega \in \Omega \mid \mathbf{X}(\omega) \in I\}),$$

definované pro libovolnou borelovskou množinu I v \mathbb{R}^n (speciálně pro libovolné sjednocení spočetně mnoha n -rozměrných intervalů) a určující **rozdělení náhodného vektoru** \mathbf{X} .

Úspornější reprezentace: Stačí intervaly tvaru $I_k = (-\infty, t_k)$, $t_k \in \mathbb{R}$,

$$\begin{aligned} P[X_1 \in (-\infty, t_1), \dots, X_n \in (-\infty, t_n)] &= P[X_1 \leq t_1, \dots, X_n \leq t_n] = \\ &= P_{\mathbf{X}}((-\infty, t_1) \times \dots \times (-\infty, t_n)) = \\ &= F_{\mathbf{X}}(t_1, \dots, t_n). \end{aligned}$$

$F_{\mathbf{X}}: \mathbb{R}^n \rightarrow \langle 0, 1 \rangle$ je **distribuční funkce** náhodného vektoru \mathbf{X} . Je

- neklesající (ve všech proměnných),
- zprava spojitá (ve všech proměnných),

- $\lim_{t_1 \rightarrow \infty, \dots, t_n \rightarrow \infty} F_{\mathbf{X}}(t_1, \dots, t_n) = 1,$
- $\forall k \in \{1, \dots, n\} \quad \forall t_1, \dots, t_{k-1}, t_{k+1}, \dots, t_n : \lim_{t_k \rightarrow -\infty} F_{\mathbf{X}}(t_1, \dots, t_n) = 0.$

Věta: Tyto podmínky jsou **nutné, nikoli postačující**.

Nestačí znát **marginální** rozdělení náhodných veličin X_1, \dots, X_n , neboť ta neobsahují informace o závislosti.

4.3 Nezávislost náhodných veličin

Náhodné veličiny X_1, X_2 jsou **nezávislé**, pokud pro všechny intervaly I_1, I_2 jsou jevy $X_1 \in I_1, X_2 \in I_2$ nezávislé, tj.

$$P[X_1 \in I_1, X_2 \in I_2] = P[X_1 \in I_1] \cdot P[X_2 \in I_2].$$

Stačí se omezit na intervaly tvaru $(-\infty, t)$, tj.

$$P[X_1 \leq t_1, X_2 \leq t_2] = P[X_1 \leq t_1] \cdot P[X_2 \leq t_2],$$

neboli

$$F_{X_1, X_2}(t_1, t_2) = F_{X_1}(t_1) \cdot F_{X_2}(t_2)$$

pro všechna $t_1, t_2 \in \mathbb{R}$.

Náhodné veličiny X_1, \dots, X_n jsou **nezávislé**, pokud pro libovolné intervaly I_1, \dots, I_n platí

$$P[X_1 \in I_1, \dots, X_n \in I_n] = \prod_{i=1}^n P[X_i \in I_i].$$

Na rozdíl od definice nezávislosti více než 2 jevů, zde není třeba požadovat nezávislost pro libovolnou podmnožinu náhodných veličin X_1, \dots, X_n . Ta vyplývá z toho, že libovolnou náhodnou veličinu X_i lze „vynechat“ tak, že zvolíme příslušný interval $I_i = \mathbb{R}$. Pak $P[X_i \in I_i] = 1$ a v součinu se tento činitel neprojeví.

Ekvivalentně stačí požadovat

$$P[X_1 \leq t_1, \dots, X_n \leq t_n] = \prod_{i=1}^n P[X_i \leq t_i]$$

pro všechna $t_1, \dots, t_n \in \mathbb{R}$, což pro sdruženou distribuční funkci **nezávislých** náhodných veličin znamená

$$F_{\mathbf{X}}(t_1, \dots, t_n) = \prod_{k=1}^n F_{X_k}(t_k).$$

Náhodné veličiny X_1, \dots, X_n jsou **po dvou nezávislé**, pokud každé dvě (různé) z nich jsou nezávislé. To je slabší podmínka než **nezávislost** veličin X_1, \dots, X_n .

4.4 Obecnější náhodné veličiny

Komplexní náhodná veličina je náhodný vektor se dvěma složkami interpretovanými jako reálná a imaginární část.

Někdy připouštíme i „náhodné veličiny“, jejichž hodnoty jsou jiné než numerické. Mohou to být např. náhodné množiny. Jindy nabývají konečně mnoha hodnot, kterým ponecháme jejich přirozené označení, např. „rub“, „líc“, „kámen“, „nůžky“, „papír“ apod.

Na těchto hodnotách nemusí být definovaná žádná aritmetika ani uspořádání.

Mohli bychom všechny hodnoty očíslovat, ale není žádný důvod, proč bychom to měli udělat právě určitým způsobem (který by ovlivnil následné numerické výpočty).

(Příklad: Číslování politických stran ve volbách.)

4.5 Směs náhodných veličin

Příklad: Náhodné veličiny U, V jsou výsledky studenta při odpovědích na dvě zkuškové otázky. Učitel náhodně vybere s pravděpodobností c první otázku, s pravděpodobností $1 - c$ druhou; podle odpovědi na vybranou otázku udělí známku. Jaké rozdělení má výsledná známka X ?

Matematický model vyžaduje vytvoření odpovídajícího pravděpodobnostního prostoru pro tento pokus.

Nechť U , resp. V je náhodná veličina na pravděpodobnostním prostoru $(\Omega_1, \mathcal{A}_1, P_1)$, resp. $(\Omega_2, \mathcal{A}_2, P_2)$, přičemž $\Omega_1 \cap \Omega_2 = \emptyset$.

Nechť $c \in \langle 0, 1 \rangle$.

Definujeme nový pravděpodobnostní prostor (Ω, \mathcal{A}, P) , kde

$\Omega = \Omega_1 \cup \Omega_2$, $\mathcal{A} = \{A_1 \cup A_2 \mid A_1 \in \mathcal{A}_1, A_2 \in \mathcal{A}_2\}$,

$P(A_1 \cup A_2) = cP_1(A_1) + (1 - c)P_2(A_2)$ pro $A_1 \in \mathcal{A}_1, A_2 \in \mathcal{A}_2$.

Definujeme funkci $X: \Omega \rightarrow \mathbb{R}$:

$$X(\omega) = \begin{cases} U(\omega) & \text{pro } \omega \in \Omega_1, \\ V(\omega) & \text{pro } \omega \in \Omega_2. \end{cases}$$

X je náhodná veličina na (Ω, \mathcal{A}, P) .

X nazýváme **směs náhodných veličin** U, V s **koeficientem** c (angl. *mixture*), značíme $\text{Mix}_c(U, V)$. Má pravděpodobnostní míru

$$P_X = cP_U + (1 - c)P_V$$

a distribuční funkci

$$F_X = cF_U + (1 - c)F_V.$$

Podobně definujeme obecněji **směs náhodných veličin** U_1, \dots, U_n s **koeficienty** $c_1, \dots, c_n \in \langle 0, 1 \rangle$, $\sum_{i=1}^n c_i = 1$, značíme $\text{Mix}_{(c_1, \dots, c_n)}(U_1, \dots, U_n) = \text{Mix}_c(U_1, \dots, U_n)$, kde $\mathbf{c} = (c_1, \dots, c_n)$.

Má pravděpodobnostní míru $\sum_{i=1}^n c_i P_{U_i}$ a distribuční funkci $\sum_{i=1}^n c_i F_{U_i}$. (Lze zobecnit i na spočetně mnoho náhodných veličin.)

Podíl jednotlivých složek je určen vektorem koeficientů $\mathbf{c} = (c_1, \dots, c_n)$. Jejich počet je stejný jako počet náhodných veličin ve směsi. Jelikož $c_n = 1 - \sum_{i=1}^{n-1} c_i$, poslední koeficient někdy vynecháváme.

Speciálně pro dvě náhodné veličiny $\text{Mix}_{(c, 1-c)}(U, V) = \text{Mix}_c(U, V)$ (kde c je číslo, nikoli vektor).

Příklad: Směsí reálných čísel r_1, \dots, r_n s koeficienty c_1, \dots, c_n je náhodná veličina $X = \text{Mix}_{(c_1, \dots, c_n)}(r_1, \dots, r_n)$,

$$P_X(I) = P[X \in I] = \sum_{i: r_i \in I} c_i, \quad F_X(t) = \sum_{i: r_i \leq t} c_i.$$

Lze ji popsat též **pravděpodobnostní funkcí** $p_X: \mathbb{R} \rightarrow \langle 0, 1 \rangle$,

$$p_X(t) = P_X(\{t\}) = P[X = t] = \begin{cases} c_i & \text{pro } t = r_i, \\ 0 & \text{jinak} \end{cases}$$

(pokud jsou r_1, \dots, r_n navzájem různá). Možno zobecnit i na spočetně mnoho reálných čísel.

4.6 Druhy náhodných veličin

- Diskrétní:** (z předchozího příkladu) Existuje spočetná množina O_X , pro kterou $P_X(\mathbb{R} \setminus O_X) = P[X \notin O_X] = 0$. Nejmenší taková množina (pokud existuje) je $\Omega_X = \{t \in \mathbb{R} : P_X(\{t\}) \neq 0\} = \{t \in \mathbb{R} : P[X = t] \neq 0\}$.

Diskrétní náhodnou veličinu popisuje **pravděpodobnostní funkce** $p_X(t) = P_X(\{t\}) = P[X = t]$.

Splňuje $\sum_{t \in \mathbb{R}} p_X(t) = 1$.

- Spojité:** Má spojitou distribuční funkci.
- Smíšená:** Směs předchozích dvou případů;
 $\Omega_X \neq \emptyset, P_X(\mathbb{R} \setminus \Omega_X) = P[X \notin \Omega_X] \neq 0$.

4.7 Popis spojitě náhodné veličiny

Náhodná veličina X je **absolutně spojitá**, jestliže existuje nezáporná funkce $f_X: \mathbb{R} \rightarrow \langle 0, \infty \rangle$ (**hustota** náhodné veličiny X) taková, že

$$F_X(t) = \int_{-\infty}^t f_X(u) du.$$

Hustota splňuje $\int_{-\infty}^{\infty} f_X(u) du = 1$.

Není určena jednoznačně, ale dvě hustoty f_X, g_X téže náhodné veličiny splňují $\int_I (f_X(x) - g_X(x)) dx = 0$ pro všechny intervaly I .

Lze volit $f_X(t) = \frac{dF_X(t)}{dt}$, pokud derivace existuje.

$P_X(\{t\}) = 0$ pro všechna t .

Některé **spojité** náhodné veličiny nejsou **absolutně spojité**; mají spojitou distribuční funkci, kterou nelze vyjádřit jako integrál. Tyto případy dále neuvažujeme.

4.8 Popis smíšené náhodné veličiny

Náhodnou veličinu X se smíšeným rozdělením nelze popsat ani pravděpodobnostní funkcí (existuje, ale neurčuje celé rozdělení) ani hustotou (neexistuje, nevychází konečná), ale lze ji **jednoznačně** vyjádřit ve tvaru $X = \text{Mix}_c(U, V)$, kde U je diskrétní, V je spojitá a $c \in (0, 1)$:

$$c = P_X(\Omega_X) = P_X(\{t \in \mathbb{R} : P_X(\{t\}) \neq 0\}),$$

$$c P_U(\{t\}) + (1 - c) \underbrace{P_V(\{t\})}_0 = c P_U(\{t\}) = P_X(\{t\}),$$

$$p_U(t) = P_U(\{t\}) = \frac{P_X(\{t\})}{c},$$

$$\Omega_U = \Omega_X,$$

$$c P_U(I) + (1 - c) P_V(I) = P_X(I),$$

$$P_V(I) = \frac{P_X(I) - c P_U(I)}{1 - c},$$

$$F_V(t) = \frac{F_X(t) - c F_U(t)}{1 - c}.$$

Alternativa bez použití pravděpodobnostní míry:

$$p_X(t) = P[X = t] = \lim_{u \rightarrow t^+} F_X(t) - \lim_{u \rightarrow t^-} F_X(t),$$

$$c = \sum_{t \in \mathbb{R}} p_X(t),$$

$$c p_U(t) = p_X(t),$$

$$p_U(t) = \frac{p_X(t)}{c},$$

$$c F_U(t) + (1 - c) F_V(t) = F_X(t),$$

$$F_V(t) = \frac{F_X(t) - c F_U(t)}{1 - c}.$$

(Lze ještě pokračovat rozkladem diskrétní části na směs Diracových rozdělení.)

4.9 Kvantilová funkce náhodné veličiny

Příklad 1. Pokud absolvent školy říká, že patří mezi 5% nejlepších, pak tvrdí, že distribuční funkce prospěchu (náhodně vybraného absolventa) má u jeho prospěchu hodnotu nejvýše 0.05. (Předpokládáme, že lepšímu prospěchu odpovídá nižší průměr známek.)

Neostrá nerovnost v definici znamená, že hodnota distribuční funkce udává podíl těch absolventů, kteří měli lepší nebo stejný prospěch.

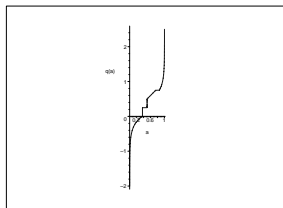
Obráceně se lze ptát, jaký prospěch je potřeba k tomu, aby se absolvent dostal mezi 5% nejlepších.

Pro $\alpha \in (0, 1)$ hledáme $t \in \mathbb{R}$ takové, že $F_X(t) = \alpha$. Máme však zaručeno pouze, že

$$\exists t \in \mathbb{R} : P[X < t] \leq \alpha \leq P[X \leq t] = F_X(t).$$

Všechna taková t tvoří omezený interval a vezmeme z něj (obvykle) střed, přesněji tedy

$$q_X(\alpha) = \frac{1}{2} (\sup \{t \in \mathbb{R} \mid P[X < t] \leq \alpha\} + \inf \{t \in \mathbb{R} \mid P[X \leq t] \geq \alpha\}).$$



Číslo $q_X(\alpha)$ se nazývá α -**kvantil** náhodné veličiny X a funkce $q_X: (0, 1) \rightarrow \mathbb{R}$ je **kvantilová funkce** náhodné veličiny X . Speciálně $q_X(\frac{1}{2})$ je **medián**, další kvantily mají také svá jména – **tercil**, **kvartil** (**dolní** $q_X(\frac{1}{4})$, **horní** $q_X(\frac{3}{4})$) ... **decil** ... **centil** neboli **percentil**

Vlastnosti kvantilové funkce:

- neklesající,
- $q_X(\alpha) = \frac{1}{2} (q_X(\alpha-) + q_X(\alpha+))$.

Věta: Tyto podmínky jsou **nutné** i **postačující**.

Obrácený převod:

$$F_X(t) = \inf\{\alpha \in (0, 1) \mid q_X(\alpha) > t\} = \sup\{\alpha \in (0, 1) \mid q_X(\alpha) \leq t\}.$$

Funkce F_X, q_X jsou navzájem inverzní tam, kde jsou spojité a rostoucí (tyto podmínky stačí ověřit pro jednu z nich).

4.10 Jak reprezentovat náhodnou veličinu v počítači

1. **Diskrétní:** Nabývá-li pouze konečného počtu hodnot t_k , $k = 1, \dots, n$, stačí k reprezentaci tyto hodnoty a jejich pravděpodobnosti $p_X(t_k) = P_X(\{t_k\}) = P[X = t_k]$, čímž je plně popsána pravděpodobnostní funkce $2n$ čísl (až na nepřesnost zobrazení reálných čísel v počítači).

Pokud diskrétní náhodná veličina nabývá (spočetně) nekonečně mnoha hodnot, musíme některé vynechat, zejména ty, které jsou málo pravděpodobné. Pro každé $\varepsilon > 0$ lze vybrat konečně mnoho hodnot t_k , $k = 1, \dots, n$, tak, že $P_X(\mathbb{R} \setminus \{t_1, \dots, t_n\}) = P[X \notin \{t_1, \dots, t_n\}] \leq \varepsilon$. Zbývá však problém, jakou hodnotu přiřadit zbývajícím (byť málo pravděpodobným) případům.

2. **(Absolutně) spojitá:** Hustotu můžeme přibližně popsat hodnotami $f(t_k)$ v „dostatečně mnoha“ bodech t_k , $k = 1, \dots, n$, ale jen za předpokladu, že je „dostatečně hladká“. Zajímají nás z ní spíše integrály typu

$$F_X(t_{k+1}) - F_X(t_k) = \int_{t_k}^{t_{k+1}} f_X(u) du,$$

z nichž lze přibližně zkonstruovat distribuční funkci. Můžeme pro reprezentaci použít přímo hodnoty distribuční funkce $F_X(t_k)$. Tam, kde je hustota velká, potřebujeme volit body hustě.

Můžeme volit body t_k , $k = 1, \dots, n$, tak, aby přírůstky $F_X(t_{k+1}) - F_X(t_k)$ měly zvolenou velikost. Zvolíme tedy $\alpha_k \in (0, 1)$, $k = 1, \dots, n$, a k nim najdeme čísla $t_k = q_X(\alpha_k)$.

Paměťová náročnost je velká, závisí na jemnosti škály hodnot náhodné veličiny, resp. její distribuční funkce.

Často je rozdělení známého typu a stačí doplnit několik parametrů, aby bylo plně určeno.

Mnohé obecnější případy se snažíme vyjádřit alespoň jako směsi náhodných veličin s rozděleními známého typu, abychom vystačili s konečně mnoha parametry.

- Smíšená:** Jako u spojitě náhodné veličiny. Tento popis je však pro diskrétní část zbytečně nepřesný.

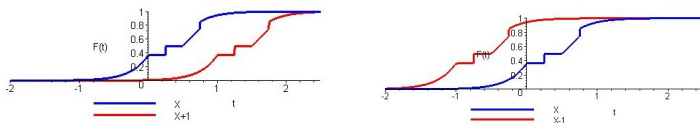
Můžeme použít rozklad na diskrétní a spojitou část.

4.11 Operace s náhodnými veličinami

Zde $I, J \subseteq \mathbb{R}$ jsou intervaly nebo početná sjednocení intervalů.

Přičtení konstanty r odpovídá posunutí ve směru vodorovné osy:

$$\begin{aligned} P_{X+r}(I+r) &= P_X(I), & P_{X+r}(J) &= P_X(J-r), \\ F_{X+r}(t+r) &= F_X(t), & F_{X+r}(u) &= F_X(u-r), \\ q_{X+r}(\alpha) &= q_X(\alpha) + r. \end{aligned}$$

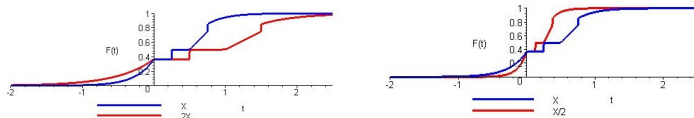


Vynásobení nenulovou konstantou r odpovídá podobnost ve směru vodorovné osy:

$$P_{rX}(rI) = P_X(I), \quad P_{rX}(J) = P_X\left(\frac{J}{r}\right).$$

Pro distribuční funkci musíme rozlišit případy:

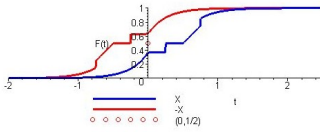
- $r > 0$: $F_{rX}(rt) = F_X(t)$, $F_{rX}(u) = F_X\left(\frac{u}{r}\right)$, $q_{rX}(\alpha) = r q_X(\alpha)$,



- $r = -1$: $F_{-X}(-t) = P_{-X}((-\infty, -t)) = P_X((t, \infty)) = 1 - P_X((-\infty, t))$, **v bodech spojitosti** distribuční funkce $F_{-X}(-t) = 1 - P_X((-\infty, t)) = 1 - P[X < t] = 1 - P[X \leq t] = 1 - P_X((-\infty, t)) = 1 - F_X(t)$,

$F_{-X}(u) = 1 - F_X(-u)$, **v bodech nespojitosti** limita zprava (středová symetrie grafu podle bodu $(0, \frac{1}{2})$ s opravou na spojitost zprava),

$$q_{-X}(\alpha) = -q_X(1 - \alpha),$$



- $r < 0$: kombinace předchozích případů.

Zobrazení spojitou rostoucí funkcí h :

$$F_{h(X)}(h(I)) = P_X(I), \quad F_{h(X)}(h(t)) = F_X(t), \quad F_{h(X)}(u) = F_X(h^{-1}(u)),$$

$q_{h(X)}(\alpha) = h(q_X(\alpha))$ v bodech spojitosti kvantilové funkce.

Zobrazení neklesající, zleva spojitou funkcí h :

$$F_{h(X)}(u) = \sup\{F_X(t) \mid h(t) \leq u\}.$$

Zobrazení po částech monotonní, zleva spojitou funkcí h :

Můžeme vyjádřit $h = h_+ - h_-$, kde h_+, h_- jsou neklesající.

X vyjádříme jako směs $X = \text{Mix}_c(U, V)$, kde U nabývá pouze hodnot, v nichž je h neklesající, V pouze hodnot, v nichž je h nerostoucí. Výsledek dostaneme jako směs dvou náhodných veličin, vzniklých zobrazením funkcemi h_+, h_- . Funkci h lze aplikovat na směs „po složkách“, tj. $h(\text{Mix}_c(U, V)) = \text{Mix}_c(h(U), h(V))$.

Součet náhodných veličin není jednoznačně určen, jedině za předpokladu **nezávislosti**.

Ani pak není vztah jednoduchý.

Směs náhodných veličin viz výše. Na rozdíl od součtu je plně určena (marginálními) rozděleními vstupních náhodných veličin a koeficienty směsi.

4.12 Jak realizovat náhodnou veličinu na počítači

1. Vytvoříme náhodný (nebo pseudonáhodný) generátor náhodné veličiny X s rovnoměrným rozdělením na $\langle 0, 1 \rangle$.
2. Náhodná veličina $q_Y(X)$ má stejné rozdělení jako Y . (Stačí tedy na každou realizaci náhodné veličiny X aplikovat funkci q_Y .)

Všechna rozdělení **spojitých** náhodných veličin jsou stejná až na (nelineární) změnu měřítka.

4.13 Střední hodnota

Značení: E . nebo μ .

Je definována zvláště pro

- **diskrétní** náhodnou veličinu U :

$$EU = \mu_U = \sum_{t \in \mathbb{R}} t \cdot p_U(t) = \sum_{t \in \Omega_U} t \cdot p_U(t),$$

- **spojitou** náhodnou veličinu V :

$$EV = \mu_V = \int_{-\infty}^{\infty} t \cdot f_V(t) dt,$$

- **směs** náhodných veličin $X = \text{Mix}_c(U, V)$, kde U je diskrétní, V je spojitá:

$$EX = cEU + (1 - c)EV .$$

(To **není** linearita střední hodnoty!)

Lze vyjít z definice pro diskrétní náhodnou veličinu a ostatní případy dostat jako limitu pro aproximaci jiných rozdělení diskrétním.

Všechny tři případy pokrývá univerzální vzorec s použitím kvantilové funkce

$$EX = \int_0^1 q_X(\alpha) d\alpha .$$

Ten lze navíc jednoduše zobecnit na střední hodnotu jakékoli funkce náhodné veličiny:

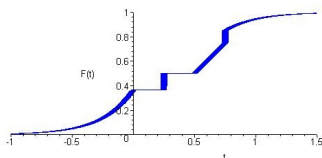
$$E(h(X)) = \int_0^1 h(q_X(\alpha)) d\alpha .$$

Speciálně pro **diskrétní** náhodnou veličinu

$$E(h(U)) = \sum_{t \in \Omega_U} h(t) \cdot p_U(t) ,$$

pro spojitou náhodnou veličinu by obdobný vzorec platil jen za omezujících předpokladů, protože spojitost náhodné veličiny se nemusí zachovávat.

Střední hodnota je vodorovnou souřadnicí těžiště grafu distribuční funkce, jsou-li jeho elementy váženy přírůstkem distribuční funkce:



Pokud pracujeme se střední hodnotou, automaticky předpokládáme, že existuje (což není vždy splněno).

4.13.1 Vlastnosti střední hodnoty

$$\begin{aligned} Er &= r, & \text{spec.} & \quad E(EX) = EX, \\ E(X + Y) &= EX + EY, & \text{spec.} & \quad E(X + r) = EX + r, \\ E(X - Y) &= EX - EY, \\ E(rX) &= rEX, & \text{obecněji} & \quad E(rX + sY) = rEX + sEY. \end{aligned}$$

(To **je** linearita střední hodnoty.)

$$E(\text{Mix}_c(U, V)) = cEU + (1 - c)EV .$$

(To **není** linearita střední hodnoty.)

Pouze pro **nezávislé** náhodné veličiny

$$E(X \cdot Y) = EX \cdot EY .$$

4.14 Rozptyl (disperze)

Značení: σ^2 , D., var.

$$DX = E \left((X - EX)^2 \right) = E(X^2) - (EX)^2,$$

$$E(X^2) = (EX)^2 + DX. \quad (1)$$

Vlastnosti:

$$DX = \int_0^1 (q_X(\alpha) - EX)^2 d\alpha.$$

$$DX \geq 0,$$

$$D r = 0,$$

$$D(X + r) = DX,$$

$$D(r X) = r^2 DX.$$

$$\begin{aligned} D(\text{Mix}_c(U, V)) &= E(\text{Mix}_c(U, V)^2) - (E(\text{Mix}_c(U, V)))^2 \\ &= cE(U^2) + (1-c)E(V^2) - (cEU + (1-c)EV)^2 \\ &= c(DU + (EU)^2) + (1-c)(DV + (EV)^2) \\ &\quad - (c^2(EU)^2 + 2c(1-c)EUEV + (1-c)^2(EV)^2) \\ &= cDU + (1-c)DV + c(1-c)(EU)^2 \\ &\quad - 2c(1-c)EUEV + c(1-c)(EV)^2 \\ &= cDU + (1-c)DV + c(1-c)(EU - EV)^2. \end{aligned}$$

Pouze pro **nezávislé** náhodné veličiny

$$D(X + Y) = DX + DY, \quad D(X - Y) = DX + DY.$$

4.15 Směrodatná odchylka

Značení: σ .

$$\sigma_X = \sqrt{DX} = \sqrt{E \left((X - EX)^2 \right)}$$

Na rozdíl od rozptylu má **stejný fyzikální rozměr** jako původní náhodná veličina.

Vlastnosti:

$$\sigma_X = \sqrt{\int_0^1 (q_X(\alpha) - EX)^2 d\alpha}.$$

$$\begin{aligned}\sigma_X &\geq 0, \\ \sigma_r &= 0, \\ \sigma_{X+r} &= \sigma_X, \\ \sigma_{rX} &= |r| \sigma_X.\end{aligned}$$

Pouze pro **nezávislé** náhodné veličiny

$$\sigma_{X+Y} = \sqrt{DX + DY} = \sqrt{\sigma_X^2 + \sigma_Y^2}.$$

4.16 Obecné a centrální momenty

$k \in \mathbb{N}$
 k -tý **obecný moment** (značení *nezavádíme*): $E(X^k)$, speciálně:
 pro $k = 1$: EX ,
 pro $k = 2$: $E(X^2) = (EX)^2 + DX$.
 Alternativní značení: m_k, μ'_k .

k -tý **centrální moment** (značení *nezavádíme*): $E((X - EX)^k)$, speciálně:
 pro $k = 1$: 0 ,
 pro $k = 2$: DX .
 Alternativní značení: μ_k .

Pomocí kvantilové funkce:

$$\begin{aligned}E(X^k) &= \int_0^1 (q_X(\alpha))^k d\alpha. \\ E((X - EX)^k) &= \int_0^1 (q_X(\alpha) - EX)^k d\alpha.\end{aligned}$$

4.17 Normovaná náhodná veličina

je taková, která má nulovou střední hodnotu a jednotkový rozptyl:

$$\text{norm } X = \frac{X - EX}{\sigma_X}$$

(pokud má vzorec smysl). Zpětná transformace je

$$X = EX + \sigma_X \text{ norm } X. \quad (2)$$

4.18 Základní typy diskrétních rozdělení

4.18.1 Diracovo

Je jediný možný výsledek $r \in \mathbb{R}$.

$$p_X(r) = 1, \quad EX = r, \quad DX = 0.$$

Všechna diskrétní rozdělení jsou směsí Diracových rozdělení.

4.18.2 Rovnoměrné

Je m možných výsledků stejně pravděpodobných.
Speciálně pro obor hodnot $\{1, 2, \dots, m\}$ dostáváme

$$p_X(k) = \frac{1}{m}, \quad k \in \{1, 2, \dots, m\},$$
$$EX = \frac{m+1}{2}, \quad DX = \frac{1}{12} (m+1)(m-1).$$

4.18.3 Alternativní (Bernoulliovo)

Jsou 2 možné výsledky. (Směs dvou Diracových rozdělení.)
Pokud výsledky jsou 0, 1, kde 1 má pravděpodobnost $q \in (0, 1)$, dostáváme

$$p_X(1) = q, \quad p_X(0) = 1 - q,$$
$$EX = q, \quad DX = q(1 - q).$$

4.18.4 Binomické $Bi(m, q)$

Počet úspěchů z m nezávislých pokusů, je-li v každém stejná pravděpodobnost úspěchu $q \in (0, 1)$. (Součet m nezávislých alternativních rozdělení.)

$$p_X(k) = \binom{m}{k} q^k (1 - q)^{m-k}, \quad k \in \{0, 1, 2, \dots, m\},$$
$$EX = m q, \quad DX = m q (1 - q).$$

Výpočetní složitost výpočtu $p_X(k)$ je $O(k)$, celého rozdělení $O(m^2)$.

4.18.5 Poissonovo $Po(\lambda)$

Limitní případ binomického rozdělení pro $m \rightarrow \infty$ při konstantním $m q = \lambda > 0$ (tedy $q \rightarrow 0$).

$$p_X(k) = \frac{\lambda^k}{k!} e^{-\lambda}, \quad k \in \{0, 1, 2, \dots\}.$$

Jednotlivé pravděpodobnosti se počítají snáze než u binomického rozdělení (ovšem všechny nevypočítáme, protože jich je nekonečně mnoho).

$$EX = \lambda, \quad DX = \lambda.$$

*„Střední hodnota se rovná rozptylu;“ jedná se **vždy o bezrozměrné celočíselné** náhodné veličiny (počet výskytů).*

Poissonovo rozdělení jako limitní případ binomického Pro $m \rightarrow \infty$ při konstantním $m q = \lambda$, tj. $q = \frac{\lambda}{m}$:

$$\begin{aligned} p_X(k) &= \binom{m}{k} q^k (1-q)^{m-k} = \\ &= \frac{m(m-1)\dots(m-(k-1))}{k!} \left(\frac{\lambda}{m}\right)^k \left(1 - \frac{\lambda}{m}\right)^{m-k} = \\ &= \frac{\lambda^k}{k!} \underbrace{1 \left(1 - \frac{1}{m}\right) \dots \left(1 - \frac{k-1}{m}\right)}_{\rightarrow 1} \underbrace{\left(1 - \frac{\lambda}{m}\right)^{-k}}_{\rightarrow 1} \underbrace{\left(1 - \frac{\lambda}{m}\right)^m}_{\rightarrow e^{-\lambda}} \rightarrow \\ &\rightarrow \frac{\lambda^k}{k!} e^{-\lambda}. \end{aligned}$$

4.18.6 Geometrické

Počet úspěchů do prvního neúspěchu, je-li v každém pokusu stejná pravděpodobnost úspěchu $q \in (0, 1)$.

$$\begin{aligned} p_X(k) &= q^k (1-q), \quad k \in \{0, 1, 2, \dots\}, \\ EX &= \frac{q}{1-q}, \quad DX = \frac{q}{(1-q)^2}. \end{aligned}$$

4.18.7 Hypergeometrické

Počet výskytů v m vzorcích, vybraných z M objektů, v nichž je K výskytů ($1 \leq m \leq K \leq M$).

$$\begin{aligned} p_X(k) &= \frac{\binom{K}{k} \binom{M-K}{m-k}}{\binom{M}{m}}, \quad k \in \{0, 1, 2, \dots, m\}, \\ EX &= \frac{mK}{M}, \quad DX = \frac{mK(M-K)(M-m)}{M^2(M-1)}. \end{aligned}$$

Výpočetní složitost výpočtu $p_X(k)$ je $O(m)$, celého rozdělení $O(m^2)$.

Binomické rozdělení jako limitní případ hypergeometrického

Lemma: Pro $m, M \in \mathbb{N}$, $m < M$, je

$$\lim_{M \rightarrow \infty} \binom{M}{m} \frac{m!}{M^m} = 1.$$

Důkaz:

$$\binom{M}{m} \frac{m!}{M^m} = \frac{M(M-1)\dots(M-(m-1))}{M^m} = 1 \left(1 - \frac{1}{M}\right) \dots \left(1 - \frac{m-1}{M}\right) \rightarrow 1.$$

Důsledek: Pro $M \gg m$ můžeme $\binom{M}{m}$ počítat přibližně jako $\frac{M^m}{m!}$.

Hypergeometrické rozdělení pro $M \rightarrow \infty$ při konstantním $\frac{K}{M} = q$, tj. $\frac{M-K}{M} = 1 - q$ (s využitím předchozího lemmatu):

$$\begin{aligned} p_X(k) &= \frac{\binom{K}{k} \binom{M-K}{m-k}}{\binom{M}{m}} \rightarrow \frac{\frac{K^k}{k!} \cdot \frac{(M-K)^{m-k}}{(m-k)!}}{\frac{M^m}{m!}} = \\ &= \frac{m!}{k! (m-k)!} \cdot \frac{K^k}{M^k} \cdot \frac{(M-K)^{m-k}}{M^{m-k}} = \binom{m}{k} q^k (1-q)^{m-k}. \end{aligned}$$

4.19 Základní typy spojitých rozdělení

4.19.1 Rovnoměrné $R(a, b)$

$$\begin{aligned} f_X(t) &= \begin{cases} \frac{1}{b-a} & \text{pro } t \in \langle a, b \rangle, \\ 0 & \text{jinak,} \end{cases} \\ F_X(u) &= \begin{cases} \frac{u-a}{b-a} & \text{pro } u \in \langle a, b \rangle, \\ 0 & \text{pro } u < a, \\ 1 & \text{pro } u > b, \end{cases} \\ q_X(\alpha) &= a + (b-a) \alpha, \\ EX &= \frac{a+b}{2}, \quad DX = \frac{1}{12} (b-a)^2. \end{aligned}$$

4.19.2 Normální (Gaussovo) $N(\mu, \sigma^2)$

A. Normované $N(0, 1)$:

$$\varphi(t) = f_{N(0,1)}(t) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right)$$

Distribuční funkce je transcendentní (Gaussův integrál) Φ ,

$$\Phi(u) = F_{N(0,1)}(u) = \int_{-\infty}^u \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right) dt,$$

kvantilová funkce Φ^{-1} je inverzní k Φ .

B. Obecné $N(\mu, \sigma^2)$:

$$f_{N(\mu, \sigma^2)}(t) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(t-\mu)^2}{2\sigma^2}\right), \quad EX = \mu, \quad DX = \sigma^2.$$

4.19.3 Logaritmickonormální $LN(\mu, \sigma^2)$

je rozdělení náhodné veličiny $X = \exp(Y)$, kde Y má $N(\mu, \sigma^2)$

$$\begin{aligned} f_X(u) &= \begin{cases} \frac{1}{u \sigma \sqrt{2\pi}} \exp\left(-\frac{(\ln u - \mu)^2}{2\sigma^2}\right) = \frac{f_{N(\mu, \sigma^2)}(\ln u)}{u} & \text{pro } u > 0, \\ 0 & \text{jinak,} \end{cases} \\ F_X(u) &= \begin{cases} F_{N(\mu, \sigma^2)}(\ln u) & \text{pro } u > 0, \\ 0 & \text{jinak,} \end{cases} \\ EX &= \exp\left(\mu + \frac{\sigma^2}{2}\right), \quad DX = (\exp(2\mu + \sigma^2)) (\exp(\sigma^2) - 1). \end{aligned}$$

4.19.4 Exponenciální $\text{Ex}(\tau)$

Např. rozdělení času do první poruchy, jestliže (podmíněná) pravděpodobnost poruchy za časový interval $\langle t, t + \delta \rangle$ závisí jen na δ , nikoli na t :

$$\begin{aligned} f_X(t) &= \begin{cases} \frac{1}{\tau} \exp\left(-\frac{t}{\tau}\right) & \text{pro } t > 0, \\ 0 & \text{jinak,} \end{cases} \\ F_X(u) &= \begin{cases} 1 - \exp\left(-\frac{u}{\tau}\right) & \text{pro } u > 0, \\ 0 & \text{jinak,} \end{cases} \\ q_X(\alpha) &= -\tau \ln(1 - \alpha), \\ EX &= \tau, \quad DX = \tau^2, \quad \sigma_X = \tau. \end{aligned}$$

4.20 Náhodné vektory 2

Náhodný vektor $\mathbf{X} = (X_1, \dots, X_n)$ je popsáný sdruženou distribuční funkcí $F_{\mathbf{X}}: \mathbb{R}^n \rightarrow \langle 0, 1 \rangle$

$$F_{\mathbf{X}}(t_1, \dots, t_n) = P[X_1 \leq t_1, \dots, X_n \leq t_n].$$

4.20.1 Diskrétní náhodný vektor

má všechny složky diskrétní. Lze jej popsat též **sdruženou pravděpodobnostní funkcí** $p_{\mathbf{X}}: \mathbb{R}^n \rightarrow \langle 0, 1 \rangle$

$$p_{\mathbf{X}}(t_1, \dots, t_n) = P[X_1 = t_1, \dots, X_n = t_n],$$

kteřá je nenulová jen ve spočetně mnoha bodech.

Diskrétní náhodné veličiny X_1, \dots, X_n jsou **nezávislé**, právě když

$$P[X_1 = t_1, \dots, X_n = t_n] = \prod_{i=1}^n P[X_i = t_i]$$

pro všechna $t_1, \dots, t_n \in \mathbb{R}$. Ekvivalentní formulace:

$$p_{\mathbf{X}}(t_1, \dots, t_n) = \prod_{i=1}^n p_{X_i}(t_i).$$

4.20.2 Spojitý náhodný vektor

má všechny složky spojité. Lze jej popsat též **sdruženou hustotou pravděpodobnosti** což je (každá) nezáporná funkce $f_{\mathbf{X}}: \mathbb{R}^n \rightarrow \langle 0, \infty \rangle$ taková, že

$$F_{\mathbf{X}}(t_1, \dots, t_n) = \int_{-\infty}^{t_1} \dots \int_{-\infty}^{t_n} f_{\mathbf{X}}(u_1, \dots, u_n) du_1 \dots du_n,$$

pro všechna $t_1, \dots, t_n \in \mathbb{R}$. Pokud to jde, volíme

$$f_{\mathbf{X}}(u_1, \dots, u_n) = \frac{\partial}{\partial t_1} \frac{\partial}{\partial t_2} \dots \frac{\partial}{\partial t_n} F_{\mathbf{X}}(t_1, \dots, t_n) = D_1 D_2 \dots D_n F_{\mathbf{X}}(t_1, \dots, t_n)$$

Speciálně pro intervaly $\langle a_i, b_i \rangle$ dostáváme

$$\begin{aligned} P[X_1 \in \langle a_1, b_1 \rangle, \dots, X_n \in \langle a_n, b_n \rangle] &= P_{\mathbf{X}}(\langle a_1, b_1 \rangle \times \dots \times \langle a_n, b_n \rangle) \\ &= \int_{a_1}^{b_1} \dots \int_{a_n}^{b_n} f_{\mathbf{X}}(u_1, \dots, u_n) du_1 \dots du_n \end{aligned}$$

Spojité náhodné veličiny X_1, \dots, X_n jsou **nezávislé**, právě když

$$f_{\mathbf{X}}(t_1, \dots, t_n) = \prod_{i=1}^n f_{X_i}(t_i) .$$

pro skoro všechna $t_1, \dots, t_n \in \mathbb{R}$.

4.21 Číselné charakteristiky náhodného vektoru

Střední hodnota

- náhodného vektoru $\mathbf{X} = (X_1, \dots, X_n)$: $\mathbf{E}\mathbf{X} := (EX_1, \dots, EX_n)$
- komplexní náhodné veličiny: $X = \Re(X) + i\Im(X)$: $\mathbf{E}X := E\Re(X) + iE\Im(X)$
- nenumerické náhodné veličiny: nemá smysl

Rozptyl náhodného vektoru $\mathbf{X} = (X_1, \dots, X_n)$: $\mathbf{D}\mathbf{X} := (DX_1, \dots, DX_n)$

Je-li U náhodná veličina, $a, b \in \mathbb{R}$, pak $aU + b$ má charakteristiky

$$\mathbf{E}(aU + b) = a\mathbf{E}U + b, \quad \mathbf{D}(aU + b) = a^2 \mathbf{D}U .$$

Na rozdíl od jednorozměrné náhodné veličiny, střední hodnota a rozptyl náhodného vektoru nedávají dostatečnou informaci pro výpočet rozptylu jeho lineárních funkcí. Proto zavádíme další charakteristiky. Např.

$$\mathbf{E}(X + Y) = \mathbf{E}X + \mathbf{E}Y ,$$

$$\begin{aligned} \mathbf{D}(X + Y) &= \mathbf{E}\left((X + Y)^2\right) - (\mathbf{E}(X + Y))^2 \\ &= \mathbf{E}(X^2 + Y^2 + 2XY) - (\mathbf{E}X + \mathbf{E}Y)^2 \\ &= \mathbf{E}(X^2) + \mathbf{E}(Y^2) + 2\mathbf{E}(XY) - \left((\mathbf{E}X)^2 + (\mathbf{E}Y)^2 + 2\mathbf{E}X\mathbf{E}Y\right) \\ &= \underbrace{\mathbf{E}(X^2) - (\mathbf{E}X)^2}_{\mathbf{D}X} + \underbrace{\mathbf{E}(Y^2) - (\mathbf{E}Y)^2}_{\mathbf{D}Y} + 2\underbrace{(\mathbf{E}(XY) - \mathbf{E}X\mathbf{E}Y)}_{\text{cov}(X,Y)} \\ &= \mathbf{D}X + \mathbf{D}Y + 2 \text{cov}(X, Y) , \end{aligned}$$

kde $\text{cov}(X, Y) := \mathbf{E}(XY) - \mathbf{E}X\mathbf{E}Y$ je **kovariance** náhodných veličin X, Y . Ekvivalentně ji lze definovat

$$\text{cov}(X, Y) = \mathbf{E}((X - \mathbf{E}X)(Y - \mathbf{E}Y)) ,$$

neboť

$$\begin{aligned} \mathbf{E}((X - \mathbf{E}X)(Y - \mathbf{E}Y)) &= \mathbf{E}(XY - X\mathbf{E}Y - Y\mathbf{E}X + \mathbf{E}X\mathbf{E}Y) \\ &= \mathbf{E}(XY) - \mathbf{E}X\mathbf{E}Y - \underbrace{\mathbf{E}X\mathbf{E}Y + \mathbf{E}X\mathbf{E}Y}_0 . \end{aligned}$$

(První vzorec je vhodnější pro výpočet.)

Pro existenci kovariance je postačující existence rozptylů DX, DY .

Vlastnosti kovariance:

$$\text{cov}(X, X) = DX, \quad \text{cov}(Y, X) = \text{cov}(X, Y),$$

$$\text{cov}(aX + b, cY + d) = ac \text{cov}(X, Y) \quad (a, b, c, d \in \mathbb{R})$$

(srovnejte s vlastnostmi rozptylu jako speciálního případu),

$$\text{speciálně} \quad \text{cov}(X, -X) = -DX.$$

Pro **nezávislé** náhodné veličiny X, Y je $\text{cov}(X, Y) = 0$.

Použitím kovariance pro **normované** náhodné veličiny vyjde **korelace**:

$$\varrho(X, Y) = \text{cov}(\text{norm } X, \text{norm } Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = E(\text{norm } X \cdot \text{norm } Y)$$

(předpokládáme, že směrodatné odchylky ve jmenovateli jsou **nenulové**).

Speciálně $\varrho(X, X) = 1$.

Vlastnosti korelace:

$$\varrho(X, X) = 1, \quad \varrho(X, -X) = -1, \quad \varrho(X, Y) \in \langle -1, 1 \rangle,$$

$$\varrho(Y, X) = \varrho(X, Y),$$

$$\varrho(aX + b, cY + d) = \text{sign}(ac) \varrho(X, Y) \quad (a, b, c, d \in \mathbb{R}, \quad a \neq 0 \neq c)$$

(až na znaménko nezáleží na prosté lineární transformaci).

$$\text{Důsledek:} \quad \varrho(aX + b, X) = \text{sign}(a).$$

Jsou-li náhodné veličiny X, Y **nezávislé**, je $\varrho(X, Y) = 0$. Obrácená implikace však neplatí (není to postačující podmínka pro nezávislost). Náhodné veličiny X, Y splňující $\varrho(X, Y) = 0$ nazýváme **nekorelované**.

Pro náhodný vektor $\mathbf{X} = (X_1, \dots, X_n)$ je definována **kovarianční matice**

$$\begin{aligned} \Sigma_{\mathbf{X}} &= \begin{bmatrix} \text{cov}(X_1, X_1) & \text{cov}(X_1, X_2) & \cdots & \text{cov}(X_1, X_n) \\ \text{cov}(X_2, X_1) & \text{cov}(X_2, X_2) & \cdots & \text{cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_n, X_1) & \text{cov}(X_n, X_2) & \cdots & \text{cov}(X_n, X_n) \end{bmatrix} \\ &= \begin{bmatrix} DX_1 & \text{cov}(X_1, X_2) & \cdots & \text{cov}(X_1, X_n) \\ \text{cov}(X_1, X_2) & DX_2 & \cdots & \text{cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_1, X_n) & \text{cov}(X_2, X_n) & \cdots & DX_n \end{bmatrix}. \end{aligned}$$

Je symetrická pozitivně semidefinitní, na diagonále má rozptyly.

Podobně je definována **korelační matice**

$$\varrho_{\mathbf{X}} = \begin{bmatrix} 1 & \varrho(X_1, X_2) & \cdots & \varrho(X_1, X_n) \\ \varrho(X_1, X_2) & 1 & \cdots & \varrho(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \varrho(X_1, X_n) & \varrho(X_2, X_n) & \cdots & 1 \end{bmatrix}.$$

Je symetrická pozitivně semidefinitní.

4.21.1 Vícerozměrné normální rozdělení $N(\mu, \Sigma)$

popisuje speciální případ náhodného vektoru, jehož složky mají normální rozdělení a mohou být korelované. Má hustotu

$$f_{N(\mu, \Sigma)}(\mathbf{t}) := \frac{1}{\sqrt{(2\pi)^n \det \mathbf{T}^{-1}}} \exp\left(-\frac{1}{2} (\mathbf{t} - \mu)^T \mathbf{T} (\mathbf{t} - \mu)\right),$$

kde $\mathbf{t} = (t_1, \dots, t_n) \in \mathbb{R}^n$,

$\mu = (\mu_1, \dots, \mu_n) \in \mathbb{R}^n$,

$\mathbf{T} \in \mathbb{R}^{n \times n}$ je matice, BÚNO symetrická.

Parametry rozdělení:

$\mu = (\mu_1, \dots, \mu_n) \in \mathbb{R}^n$ je střední hodnota náhodného vektoru,

$\Sigma := \mathbf{T}^{-1}$ je kovarianční matice, speciálně její hlavní diagonála

$(\Sigma_{11}, \Sigma_{22}, \dots, \Sigma_{nn}) \in \mathbb{R}^n$ je rozptyl náhodného vektoru,

marginální rozdělení i -té složky je $N(\mu_i, \Sigma_{ii})$;

pomocí těchto parametrů píšeme

$$f_{N(\mu, \Sigma)}(\mathbf{t}) := \frac{1}{\sqrt{(2\pi)^n \det \Sigma}} \exp\left(-\frac{1}{2} (\mathbf{t} - \mu)^T \Sigma^{-1} (\mathbf{t} - \mu)\right).$$

4.22 Lineární prostor náhodných veličin

(Ω, \mathcal{A}, P) pravděpodobnostní prostor,

\mathcal{L} lineární prostor všech náhodných veličin na (Ω, \mathcal{A}, P) , tj. \mathcal{A} -měřitelných funkcí $\Omega \rightarrow \mathbb{R}$,

sčítání náhodných veličin a jejich násobení reálným číslem = operace s funkcemi (bod po bodu),

\mathcal{L}_2 lineární podprostor všech náhodných veličin z \mathcal{L} , které mají rozptyl,

•: $\mathcal{L}_2 \times \mathcal{L}_2 \rightarrow \mathbb{R}$,

$$X \bullet Y := E(XY),$$

je bilineární (=lineární v obou argumentech) a komutativní operace, **skalární součin** (pokud ztotožníme náhodné veličiny X, Y , pro které $P[X \neq Y] = 0$; za prvky prostoru pak považujeme třídy ekvivalence místo jednotlivých náhodných veličin.),

$$\|X\| := \sqrt{X \bullet X} = \sqrt{E(X^2)}$$

je **norma**,

$$d(X, Y) := \|X - Y\| = \sqrt{E((X - Y)^2)}$$

je **metrika** (vzdálenost)

(bez předchozího ztotožnění pouze pseudometrika, mohla by být nulová i pro $X \neq Y$.)

\mathcal{L}_2 lze rozložit na 2 ortogonální podprostory:

\mathcal{R} = jednodimenzionální prostor všech konstatních náhodných veličin (tj. s Diracovým rozdělením)

\mathcal{N} = prostor všech náhodných veličin s nulovou střední hodnotou.

$E\mathbf{X}$ je kolmý průmět X do \mathcal{R} (pokud ztotožňujeme toto reálné číslo s příslušnou konstantní náhodnou veličinou, jinak souřadnice ve směru \mathcal{R}),

$X - E\mathbf{X}$ je kolmý průmět X do \mathcal{N} ,

norm $X = \frac{X - E\mathbf{X}}{\sigma_X}$ je jednotkový vektor ve směru kolmého průmětu X do \mathcal{N} ,

$\sigma_X = \|X - E\mathbf{X}\|$ je vzdálenost X od \mathcal{R} .

Z kolmosti vektorů $X - EX \in \mathcal{N}$, $EX \in \mathcal{R}$ a Pythagorovy věty plyne(1)

$$X \bullet X = \|X\|^2 = \|X - EX\|^2 + \|EX\|^2 ,$$

$$E(X^2) = DX + (EX)^2 .$$

4.22.1 Lineární podprostor \mathcal{N} náhodných veličin s nulovými středními hodnotami

Speciálně pro náhodné veličiny z \mathcal{N} :

$$\sigma_X^2 = X \bullet X ,$$

$$\sigma_X = \|X\| ,$$

$$\text{cov}(X, Y) = X \bullet Y ,$$

$$\varrho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{X \bullet Y}{\|X\| \|Y\|} = \cos \angle(X, Y) .$$

Důsledek: Náhodné veličiny X, Y s nulovými středními hodnotami jsou ortogonální, právě když jsou nekorelované.

Obecně v \mathcal{L}_2

$\varrho(X, Y)$ je kosinus úhlu průmětů X, Y do \mathcal{N} ,

$\text{cov}(X, Y) = X \bullet Y - EX \bullet EY$ je skalární součin průmětů X, Y do \mathcal{N} .

4.22.2 Lineární regrese

Úloha: Je dán náhodný vektor $\mathbf{X} = (X_1, \dots, X_n)$ a náhodná veličina Y .

(Předpokládáme, že všechny náhodné veličiny jsou z \mathcal{L}_2). Máme najít takové koeficienty c_1, \dots, c_n aby lineární kombinace $\sum_i c_i X_i$ byla co nejlepší aproximací náhodné veličiny Y ve smyslu kritéria

$$\left\| \sum_k c_k X_k - Y \right\| .$$

Řešení: K vektoru Y hledáme nejbližší bod v lineárním podprostoru, který je lineárním obalem vektorů X_1, \dots, X_n ; řešením je kolmý průmět. Ten je charakterizován tím, že vektor $\sum_i c_i X_i - Y$ je kolmý na X_j , $j = 1, \dots, n$,

$$\left(\sum_k c_k X_k - Y \right) \bullet X_j = 0 ,$$

$$\sum_i c_i (X_i \bullet X_j) = Y \bullet X_j .$$

To je soustava lineárních rovnic pro neznámé koeficienty c_1, \dots, c_n (soustava normálních rovnic).

Speciálně pro náhodné veličiny s nulovými středními hodnotami:

$$\sum_i c_i \text{cov}(X_i, X_j) = \text{cov}(Y, X_j) ,$$

takže matice soustavy je kovarianční matice $\Sigma_{\mathbf{X}}$.

4.23 Reprezentace náhodných vektorů v počítači

Obdobná jako u náhodných veličin, avšak s rostoucí dimenzí rychle roste paměťová náročnost. To by se nestalo, kdyby náhodné veličiny byly nezávislé; pak by stačilo znát marginální rozdělení.

Proto velkou úsporu může přinést i **podmíněná nezávislost**.

Pokud najdeme úplný systém jevů, které zajišťují podmíněnou nezávislost dvou náhodných veličin, pak můžeme jejich rozdělení popsat jako **směs** rozdělení nezávislých náhodných veličin (a tedy úsporněji).

4.24 Čebyševova nerovnost

Věta:

$$\forall \delta > 0 : P[|\text{norm } X| < \delta] \geq 1 - \frac{1}{\delta^2},$$

kde $\text{norm } X = \frac{X - \text{EX}}{\sigma_X}$ (pokud má výraz smysl).

Důkaz pomocí kvantilové funkce:

$$\underbrace{D(\text{norm } X)}_1 = E\left(\underbrace{(\text{norm } X)^2}_0\right) - \underbrace{(E(\text{norm } X))^2}_0,$$

$$1 = E\left((\text{norm } X)^2\right) = EY,$$

kde $Y = (\text{norm } X)^2$. Odhad pravděpodobnosti $\beta = P[|\text{norm } X| < \delta] = P[Y < \delta^2] = F_Y(\delta^2 -)$:

$$1 = EY = \int_0^1 q_Y(\alpha) \, d\alpha = \int_0^\beta \underbrace{q_Y(\alpha)}_{\geq 0} \, d\alpha + \int_\beta^1 \underbrace{q_Y(\alpha)}_{\geq \delta^2} \, d\alpha \geq (1 - \beta) \delta^2,$$

$$\beta \geq 1 - \frac{1}{\delta^2}.$$

Důkaz pomocí směsi: Vyjádříme $Y = (\text{norm } X)^2 = \text{Mix}_\beta(L, U)$, kde

L nabývá pouze hodnot z $(0, \delta^2)$,

U nabývá pouze hodnot z (δ^2, ∞) , takže $EU \geq \delta^2$,

$\beta = F_Y(\delta^2)$.

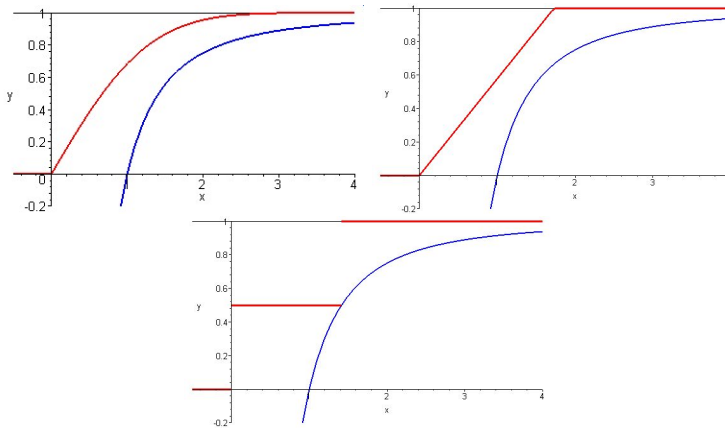
$$1 = EY = \beta \underbrace{EL}_{\geq 0} + (1 - \beta) \underbrace{EU}_{\geq \delta^2} \geq (1 - \beta) \delta^2.$$

Rovnost nastává pro $U = \delta^2$, $L = 0$, tj. pro diskrétní rozdělení $\{(EX - \delta \sigma_X, \frac{1-\beta}{2}), (EX, \beta), (EX + \delta \sigma_X, \frac{1-\beta}{2})\}$.

Ekvivalentní tvary ($\varepsilon = \delta \sigma_X$):

$$\forall \delta > 0 : P\left[\left|\frac{X - \text{EX}}{\sigma_X}\right| \geq \delta\right] \leq \frac{1}{\delta^2},$$

$$\forall \varepsilon > 0 : P[|X - \text{EX}| \geq \varepsilon] \leq \frac{\sigma_X^2}{\varepsilon^2} = \frac{DX}{\varepsilon^2}.$$



5 Základní pojmy statistiky

5.1 K čemu potřebujeme statistiku

Zkoumání **společných** vlastností velkého počtu obdobných jevů.

Přitom nezkoumáme všechny, ale jen vybraný vzorek (kvůli ceně testů, jejich destruktivnosti apod.).

- Odhady parametrů pravděpodobnostního modelu
- Testování hypotéz

Potíže statistického výzkumu – viz [Rogalewicz].

5.2 Pojem náhodného výběru, odhady

Soubor

- **základní (=populace)**
- **výběrový**

Náhodný výběr jednoho prvku **základního souboru** (s rovnoměrným rozdělením) a stanovení určitého parametru tohoto prvku určuje rozdělení náhodné veličiny.

Opakovaným výběrem dostaneme náhodný vektor, jehož složky mají stejné rozdělení a jsou nezávislé.

Takto vytvoříme **výběrový soubor rozsahu n** , obvykle však vyloučíme vícenásobný výběr stejného prvku (*výběr bez vracení*). Jeho rozdělení se může poněkud lišit od původního. Tento rozdíl se obvykle zanedbává, neboť

1. pro velký rozsah základního souboru to není podstatné,
2. rozsah základního souboru někdy není znám,
3. výpočty se značně zjednoduší.

Přesnost odhadu je dána velikostí výběrového souboru, nikoli populace.

Náhodný výběr $\mathbf{X} = (X_1, \dots, X_n)$ je vektor náhodných veličin, které jsou **nezávislé** a mají **stejné rozdělení**.

(Vynecháváme indexy, např. F_X místo F_{X_k} .)

Provedením pokusu dostaneme **realizaci náhodného výběru**,

$\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$,

kde n je **rozsah výběru**.

Statistika je (každá) měřitelná funkce G , definovaná na náhodném výběru libovolného rozsahu. (Počítá se z náhodných veličin výběru, nikoli z parametrů rozdělení.)

„**Měřitelná**“ znamená, že pro každé $t \in \mathbb{R}$ je definována pravděpodobnost

$$P[G(X_1, \dots, X_n) \leq t] = F_{G(X_1, \dots, X_n)}(t).$$

Statistika jako funkce náhodných veličin je rovněž náhodná veličina.

Obvykle se používá jako **odhad parametrů rozdělení** (které nám zůstávají skryté).

Značení:

ϑ ... skutečný parametr (reálné číslo),

$\hat{\Theta}, \hat{\Theta}_n$... jeho odhad založený na náhodném výběru rozsahu n (náhodná veličina)

$\hat{\vartheta}, \hat{\vartheta}_n$... realizace odhadu (obvykle reálné číslo)

Žádoucí vlastnosti odhadů:

- $E\hat{\Theta}_n = \vartheta$ **nestranný** (opak: **vychýlený**)
- $\lim_{n \rightarrow \infty} E\hat{\Theta}_n = \vartheta$ **asymptoticky nestranný**
- **eficientní** = s malým rozptylem, což posuzujeme podle $E\left((\hat{\Theta}_n - \vartheta)^2\right) = D\hat{\Theta}_n + \left(E(\hat{\Theta}_n - \vartheta)\right)^2$, pro nestranný odhad se redukuje na $D\hat{\Theta}_n$
- **nejlepší nestranný** odhad je ze všech nestranných ten, který je nejvíce eficientní (mohou však existovat více eficientní vychýlené odhady)
- $\lim_{n \rightarrow \infty} E\hat{\Theta}_n = \vartheta$, $\lim_{n \rightarrow \infty} \sigma_{\hat{\Theta}_n} = 0$ **konzistentní**
- **robustní**, tj. odolný vůči šumu („i při zašuměných datech dostáváme dobrý výsledek“) – zde už přesné kritérium chybí, zato je to velmi praktická vlastnost

5.3 Výběrový průměr

z náhodného výběru $\mathbf{X} = (X_1, \dots, X_n)$ je

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n X_j$$

Alternativní značení: \bar{X}_n (pokud potřebujeme zdůraznit rozsah výběru)

Jeho realizaci značíme malým písmenem:

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j.$$

Věta:

$$\begin{aligned}E\bar{\mathbf{X}}_n &= \frac{1}{n} \sum_{j=1}^n EX = EX, \\D\bar{\mathbf{X}}_n &= \frac{1}{n^2} \sum_{j=1}^n DX = \frac{1}{n} DX, \\ \sigma_{\bar{\mathbf{X}}_n} &= \sqrt{\frac{1}{n} DX} = \frac{1}{\sqrt{n}} \sigma_X,\end{aligned}$$

pokud existují. (Zde $EX = EX_j$ atd.)

Důsledek: Výběrový průměr je nestranný konzistentní odhad střední hodnoty.

(Nezávisle na typu rozdělení.)

Cebyševova nerovnost pro $\bar{\mathbf{X}}_n$ dává

$$P[|\bar{\mathbf{X}}_n - EX| \geq \varepsilon] \leq \frac{D\bar{\mathbf{X}}_n}{\varepsilon^2} = \frac{DX}{n\varepsilon^2} \rightarrow 0 \quad \text{pro } n \rightarrow \infty.$$

To platí i za obecnějších předpokladů (X_j nemusí mít stejné rozdělení) – **slabý zákon velkých čísel**.

Lidově se hovoří o „přesném součtu nepřesných čísel“, což je chyba, neboť součet $\sum_{j=1}^n X_j$ má rozptyl $nDX \rightarrow \infty$. **Relativní** chyba součtu **klesá**, **absolutní roste**.

Rozdělení výběrového průměru může být podstatně složitější než původní, jen ve speciálních případech je jednoduchá odpověď.

Věta: Výběrový průměr z **normálního** rozdělení $N(\mu, \sigma^2)$ má normální rozdělení $N(\mu, \frac{1}{n}\sigma^2)$ a je nejlepším nestranným odhadem střední hodnoty.

Podobná věta platí i pro jiná rozdělení alespoň asymptoticky:

Centrální limitní věta: Necht' X_j , $j \in \mathbb{N}$, jsou nezávislé stejně rozdělené náhodné veličiny se střední hodnotou EX a směrodatnou odchylkou $\sigma_X \neq 0$. Pak normované náhodné veličiny

$$Y_n = \text{norm } \bar{\mathbf{X}}_n = \frac{\sqrt{n}}{\sigma_X} (\bar{\mathbf{X}}_n - EX)$$

konvergují k normovanému normálnímu rozdělení v následujícím smyslu:

$$\forall t \in \mathbb{R} : \lim_{n \rightarrow \infty} F_{Y_n}(t) = \lim_{n \rightarrow \infty} F_{\text{norm } \bar{\mathbf{X}}_n}(t) = \Phi(t).$$

5.4 Výběrový rozptyl

náhodného výběru $\mathbf{X} = (X_1, \dots, X_n)$ je statistika

$$S_{\mathbf{X}}^2 = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{\mathbf{X}}_n)^2.$$

Alternativní značení: S^2 (*Dvojka v horním indexu zde neznamena kvadrát!*)

Jeho realizaci značíme malým písmenem:

$$s_{\mathbf{x}}^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{\mathbf{x}}_n)^2.$$

Praktičtější jednorůchodový vzorec:

$$S_{\mathbf{X}}^2 = \frac{1}{n-1} \sum_{j=1}^n X_j^2 - \frac{n}{n-1} \bar{\mathbf{X}}_n^2 = \frac{1}{n-1} \sum_{j=1}^n X_j^2 - \frac{1}{n(n-1)} \left(\sum_{j=1}^n X_j \right)^2.$$

Věta:

$$ES_{\mathbf{X}}^2 = DX.$$

Důkaz: Z jednorůchodového vzorce pro $S_{\mathbf{X}}^2$ dostáváme

$$\begin{aligned} ES_{\mathbf{X}}^2 &= \frac{n}{n-1} EX^2 - \frac{n}{n-1} E\bar{\mathbf{X}}_n^2 = \frac{n}{n-1} \left(DX + (EX)^2 - D\bar{\mathbf{X}}_n - (E\bar{\mathbf{X}}_n)^2 \right) = \\ &= \frac{n}{n-1} \left(DX + (EX)^2 - \frac{1}{n} DX - (EX)^2 \right) = DX. \end{aligned}$$

Věta: Výběrový rozptyl je nestranný konzistentní odhad rozptylu (pokud původní rozdělení má rozptyl a 4. centrální moment).

Rozdělení výběrového rozptylu může být podstatně složitější.

Speciálně pro rozdělení $N(0, 1)$ a $n = 2$:

$$\bar{\mathbf{X}} = \frac{X_1 + X_2}{2}, \quad X_1 - \bar{\mathbf{X}} = -(X_2 - \bar{\mathbf{X}}) = \frac{X_1 - X_2}{2} \text{ má rozdělení } N\left(0, \frac{1}{2}\right),$$

$$S_{\mathbf{X}}^2 = (X_1 - \bar{\mathbf{X}})^2 + (X_2 - \bar{\mathbf{X}})^2 = 2 \left(\frac{X_1 - X_2}{2} \right)^2 = \left(\frac{X_1 - X_2}{\sqrt{2}} \right)^2 = U^2,$$

kde $U = \frac{X_1 - X_2}{\sqrt{2}}$ má rozdělení $N(0, 1)$. Tomu říkáme:

5.4.1 Rozdělení χ^2 s 1 stupněm volnosti

= rozdělení náhodné veličiny $V = U^2$, kde U má **normované normální** rozdělení $N(0, 1)$. Značení: $\chi^2(1)$. (*Toto rozdělení není zvykem normovat.*)

$$EV = EU^2 = \underbrace{DU}_1 + \underbrace{(EU)^2}_0 = 1,$$

$$DV = 2. \quad (\text{bez důkazu})$$

Pro $t > 0$ vychází distribuční funkce

$$\begin{aligned} F_V(t) &= P[V \leq t] = P[-\sqrt{t} \leq U \leq \sqrt{t}] = 2P[0 \leq U \leq \sqrt{t}] = \\ &= 2 \left(\Phi(\sqrt{t}) - \Phi(0) \right) = 2 \int_0^{\sqrt{t}} e^{-\frac{u^2}{2}} du, \end{aligned}$$

hustota

$$f_V(t) = F'_V(t) = \left(2\Phi(\sqrt{t}) \right)' = 2(\sqrt{t})' \Phi'(\sqrt{t}) = \frac{1}{\sqrt{t}} \varphi(\sqrt{t}) = \frac{1}{\sqrt{2\pi t}} e^{-\frac{t}{2}}.$$

Zobecnění:

5.4.2 Rozdělení χ^2 s η stupni volnosti

= rozdělení náhodné veličiny $Y = \sum_{j=1}^{\eta} V_j$, kde V_j jsou **nezávislé** náhodné veličiny s rozdělením $\chi^2(1)$

= rozdělení náhodné veličiny $Y = \sum_{j=1}^{\eta} U_j^2$, kde U_j jsou **nezávislé** náhodné veličiny s **normovaným normálním** rozdělením $N(0, 1)$.

Značení: $\chi^2(\eta)$.

$$EY = E \sum_{j=1}^{\eta} V_j = \sum_{j=1}^{\eta} \underbrace{EV_j}_1 = \eta,$$

$$DY = D \sum_{j=1}^{\eta} V_j = \sum_{j=1}^{\eta} \underbrace{DV_j}_2 = 2\eta.$$

Věta: Necht' X, Y jsou **nezávislé** náhodné veličiny s rozdělením $\chi^2(\xi)$, resp. $\chi^2(\eta)$. Pak $X + Y$ má rozdělení $\chi^2(\xi + \eta)$.

Hustota

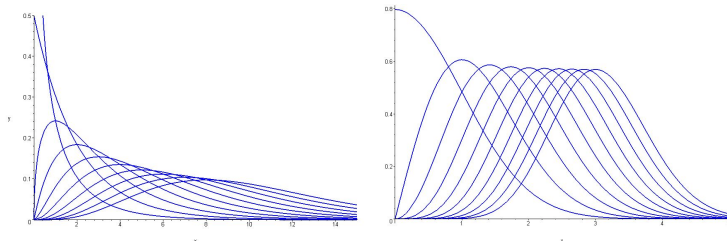
$$f_Y(y) = \begin{cases} c(\eta) y^{\frac{\eta}{2}-1} e^{-\frac{y}{2}} & \text{pro } y > 0, \\ 0 & \text{jinak,} \end{cases}$$

$$c(\eta) = \frac{1}{2^{\frac{\eta}{2}} \Gamma\left(\frac{\eta}{2}\right)},$$

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

speciálně $\Gamma(m + 1) = m!$ pro všechna $m \in \mathbb{N}$.

Speciálně pro $\eta = 2$ je $c(\eta) = 1/2$ a dostáváme exponenciální rozdělení.

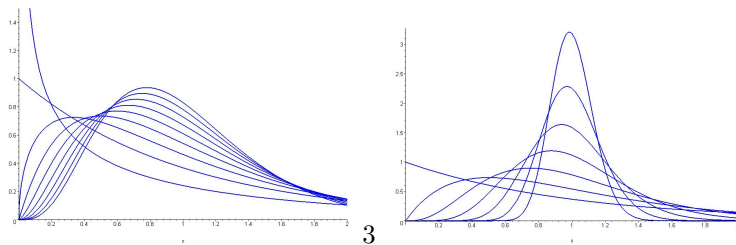


Hustoty rozdělení χ^2 s $1, 2, \dots, 10$ stupni volnosti a jeho odmocniny („vzdálenost od středu terče“).

5.4.3 Výběrový rozptyl

z **normálního** rozdělení $N(EX, DX)$ splňuje:

$$\frac{(n-1) S_X^2}{DX} \text{ má rozdělení } \chi^2(n-1).$$



Rozdělení odhadu rozptylu pomocí výběrového rozptylu $S_{\mathbf{X}}^2$ pro rozsah výběru $2, 3, \dots, 10$ a $3 = 2^1 + 1, 2^2 + 1, \dots, 2^7 + 1 = 129$.

Důsledek: Rozptyl výběrového rozptylu z normálního rozdělení $N(EX, DX)$ je

$$DS_{\mathbf{X}}^2 = \frac{2}{n-1} (DX)^2.$$

Věta: Pro náhodný výběr $\mathbf{X} = (X_1, \dots, X_n)$ z **normálního** rozdělení je $\bar{\mathbf{X}}$ nejlepší nestranný odhad střední hodnoty, $S_{\mathbf{X}}^2$ je nejlepší nestranný odhad rozptylu a statistiky $\bar{\mathbf{X}}, S_{\mathbf{X}}^2$ jsou konzistentní a **nezávislé**.

Existuje však vychýlený odhad rozptylu, který je eficientnější:

5.4.4 Alternativní odhad rozptylu

$$\widehat{DX} = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{\mathbf{X}}_n)^2 = \frac{n-1}{n} S_{\mathbf{X}}^2.$$

Věta: \widehat{DX} je vychýlený konzistentní odhad rozptylu.

Důkaz:

$$E\widehat{DX} = \frac{n-1}{n} DX \rightarrow DX,$$

\widehat{DX} má rozptyl menší než $S_{\mathbf{X}}^2$, a to v poměru $\left(\frac{n-1}{n}\right)^2$.

Eficienci nemůžeme porovnat obecně; aspoň pro **normální** rozdělení:

1. eficeence odhadu $S_{\mathbf{X}}^2$:

$$DS_{\mathbf{X}}^2 = \frac{2}{n-1} (DX)^2.$$

2. eficeence odhadu \widehat{DX} (DX je konstanta):

$$\begin{aligned} E(\widehat{DX} - DX)^2 &= D(\widehat{DX} - DX) + \left(E(\widehat{DX} - DX)\right)^2 = \\ &= D(\widehat{DX}) + \left(\frac{1}{n} DX\right)^2 = \\ &= \left(\frac{n-1}{n}\right)^2 \frac{2}{n-1} (DX)^2 + \frac{1}{n^2} (DX)^2 = \frac{2n-1}{n^2} (DX)^2, \end{aligned}$$

a protože

$$\frac{2n-1}{n^2} < \frac{2}{n} < \frac{2}{n-1},$$

je odhad \widehat{DX} více eficientní než $S_{\mathbf{X}}^2$ (který je nejlepší nestranný!).

5.5 Výběrová směrodatná odchylka

náhodného výběru $\mathbf{X} = (X_1, \dots, X_n)$ je statistika

$$S_{\mathbf{X}} = \sqrt{S_{\mathbf{X}}^2} = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X}_n)^2}.$$

Alternativní značení: S

Její realizaci značíme malým písmenem:

$$s_x = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x}_n)^2}.$$

Věta:

$$ES_{\mathbf{X}} \leq \sigma_X.$$

Rovnost obecně nenastává, takže to **není nestraný** odhad směrodatné odchylky!

Důkaz:

$$DX = ES_{\mathbf{X}}^2 = (ES_{\mathbf{X}})^2 + \underbrace{DS_{\mathbf{X}}}_{\geq 0} \geq (ES_{\mathbf{X}})^2,$$

$$\sigma_X \geq ES_{\mathbf{X}}.$$

Věta: Výběrová směrodatná odchylka je konzistentní odhad směrodatné odchylky (pokud původní rozdělení má rozptyl a 4. centrální moment).

5.6 Výběrový k -tý obecný moment

náhodného výběru $\mathbf{X} = (X_1, \dots, X_n)$ je statistika

$$M_{X^k} = \frac{1}{n} \sum_{j=1}^n X_j^k.$$

Alternativní značení: M_k

Jeho realizaci značíme malým písmenem:

$$m_{X^k} = \frac{1}{n} \sum_{j=1}^n x_j^k.$$

Věta:

$$EM_{X^k} = EX^k.$$

(Tj. je to **nestraný** odhad k -tého obecného momentu.)

Věta: Výběrový k -tý obecný moment je konzistentní odhad k -tého obecného momentu (pokud X má k -tý a $2k$ -tý obecný moment).

Důkaz:

$$DM_{X^k} = \frac{1}{n^2} n DX^k = \frac{1}{n} DX^k = \frac{1}{n} (E(X^k)^2 - (EX^k)^2) = \frac{1}{n} (EX^{2k} - (EX^k)^2).$$

5.7 Histogram a empirické rozdělení

V (nenáhodném) vektoru $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$ (získaném např. jako realizace náhodného výběru) nezáleží na pořadí složek (ale záleží na jejich opakování). Úsporněji je popsán množinou hodnot $H = \{x_1, \dots, x_n\}$ (ta má nejvýše n prvků, obvykle méně) a jejich **četnostmi** $n_t, t \in H$. Tato data obvykle znázorňujeme **tabulkou četností** nebo grafem zvaným **histogram**.

Normováním dostaneme **relativní četnosti** $r_t = \frac{n_t}{n}, t \in H$. Jelikož $\sum_{t \in H} r_t = 1$, definují relativní četnosti pravděpodobnostní funkci $p_{\text{Emp}(\mathbf{x})}(t) = r_t$ tzv. **empirického rozdělení** $\text{Emp}(\mathbf{x})$. Je to diskrétní rozdělení s nejvýše n hodnotami charakterizující vektor \mathbf{x} .

5.7.1 Vlastnosti empirického rozdělení

(*Indexem $\text{Emp}(\mathbf{x})$ označujeme parametry jakékoli náhodné veličiny, která má toto rozdělení.*)

$$\mathbb{E} \text{Emp}(\mathbf{x}) = \sum_{t \in H} t r_t = \frac{1}{n} \sum_{t \in H} t n_t = \frac{1}{n} \sum_{i=1}^n x_i = \bar{\mathbf{x}},$$

$$\mathbb{E} (\text{Emp}(\mathbf{x}))^k = \sum_{t \in H} t^k r_t = \frac{1}{n} \sum_{t \in H} t^k n_t = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

$$\begin{aligned} \text{D} \text{Emp}(\mathbf{x}) &= \sum_{t \in H} (t - \mathbb{E} \text{Emp}(\mathbf{x}))^2 r_t = \frac{1}{n} \sum_{t \in H} (t - \bar{\mathbf{x}})^2 n_t \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{\mathbf{x}})^2 = \frac{n-1}{n} s_{\mathbf{x}}^2. \end{aligned}$$

Obecné momenty empirického rozdělení se rovnají výběrovým momentům původního rozdělení. Výpočet z histogramu (z empirického rozdělení) může být jednodušší než z původní realizace náhodného výběru (pokud se opakují stejné hodnoty).

Rozptyl empirického rozdělení odpovídá odhadu $\widehat{\text{D}\bar{X}} = \frac{n-1}{n} S_{\mathbf{X}}^2$ rozptylu původního rozdělení, odlišnému od $S_{\mathbf{X}}^2$.

5.8 Výběrový medián

je medián empirického rozdělení, $q_{\text{Emp}(\mathbf{x})}(\frac{1}{2})$. Poskytuje jinou informaci než výběrový průměr, mnohdy užitečnější (mj. **robustnější** – odolnější vůči vlivu vychýlených hodnot, outliers). Navíc víme, jak se změní monotonní funkcí.

Proč se používá méně než výběrový průměr:

- Výpočetní náročnost je vyšší; seřazení hodnot má pracnost úměrnou $n \ln n$, zatímco výběrový průměr n .
- Paměťová náročnost je vyšší – potřebujeme zapamatovat všechna data, u výběrového průměru stačí 2 registry.
- Možnosti decentralizace a paralelizace výpočtu výběrového mediánu jsou velmi omezené.

5.9 Intervalové odhady

Dosud jsme skutečnou hodnotu parametru ϑ nahrazovali **bodovým odhadem** $\hat{\Theta}$ (což je náhodná veličina). Nyní místo toho hledáme **intervalový odhad**, tzv. **interval spolehlivosti** I , což je minimální interval takový, že

$$P[\vartheta \in I] \geq 1 - \alpha,$$

kde $\alpha \in (0, \frac{1}{2})$ je pravděpodobnost, že meze intervalu I budou překročeny; $1 - \alpha$ je **koeficient spolehlivosti**. Obvykle hledáme **horní**, resp. **dolní jednostranný** odhad, kdy

$$I = (-\infty, q_{\hat{\Theta}}(1 - \alpha)), \text{ resp. } I = \langle q_{\hat{\Theta}}(\alpha), \infty \rangle,$$

nebo (**symetrický**) **oboustranný** odhad,

$$I = \left\langle q_{\hat{\Theta}}\left(\frac{\alpha}{2}\right), q_{\hat{\Theta}}\left(1 - \frac{\alpha}{2}\right) \right\rangle.$$

K tomu potřebujeme znát rozdělení odhadu $\hat{\Theta}$.

5.10 Intervalové odhady parametrů **normálního** rozdělení $N(\mu, \sigma^2)$

5.10.1 Odhad střední hodnoty při **známém** rozptylu σ^2

μ odhadneme výběrovým průměrem \bar{X} s rozdělením $N\left(\mu, \frac{\sigma^2}{n}\right)$.

Normovaná náhodná veličina $\text{norm } \bar{X} = \frac{\sqrt{n}}{\sigma} (\bar{X} - \mu)$, stejně jako $-\text{norm } \bar{X} = \frac{\sqrt{n}}{\sigma} (\mu - \bar{X})$ má rozdělení $N(0, 1)$;

$$\begin{aligned} & P\left[\frac{\sqrt{n}}{\sigma} (\mu - \bar{X}) \in (-\infty, \Phi^{-1}(1 - \alpha))\right] \\ &= 1 - \alpha \\ &= P\left[\frac{\sqrt{n}}{\sigma} (\mu - \bar{X}) \leq \Phi^{-1}(1 - \alpha)\right] \\ &= P\left[\mu \leq \bar{X} + \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1 - \alpha)\right] \\ &= P\left[\mu \in \left(-\infty, \bar{X} + \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1 - \alpha)\right)\right]. \end{aligned}$$

Obdobně dostaneme i další intervalové odhady

$$\begin{aligned} & \left(-\infty, \bar{X} + \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1 - \alpha)\right), \\ & \left\langle \bar{X} - \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1 - \alpha), \infty \right\rangle, \\ & \left\langle \bar{X} - \frac{\sigma}{\sqrt{n}} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right), \bar{X} + \frac{\sigma}{\sqrt{n}} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \right\rangle, \end{aligned}$$

kde $\bar{X} - \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1 - \alpha) = \bar{X} + \frac{\sigma}{\sqrt{n}} \Phi^{-1}(\alpha)$

($\Phi^{-1}(\alpha) = -\Phi^{-1}(1 - \alpha)$ ovšem nebývá v tabulkách).

Při výpočtu nahradíme výběrový průměr \bar{X} jeho realizací \bar{x} .

5.10.2 Odhad střední hodnoty při **neznámém** rozptylu

μ odhadneme výběrovým průměrem \bar{X} s rozdělením $N\left(\mu, \frac{\sigma^2}{n}\right)$,

σ^2 odhadneme výběrovým rozptylem $S_{\mathbf{X}}^2$; $\frac{(n-1)S_{\mathbf{X}}^2}{\sigma^2}$ má rozdělení $\chi^2(n-1)$.

Testujeme analogicky náhodnou veličinu $\frac{\sqrt{n}}{S_{\mathbf{X}}}(\bar{X} - \mu)$, její rozdělení však není normální, ačkoli \bar{X} , $S_{\mathbf{X}}$ jsou nezávislé.

5.10.3 Studentovo t-rozdělení (autor: Gossett)

s η stupni volnosti je rozdělení náhodné veličiny

$$\frac{U}{\sqrt{\frac{V}{\eta}}},$$

kde U má rozdělení $N(0, 1)$,

V má rozdělení $\chi^2(\eta)$,

U, V jsou nezávislé.

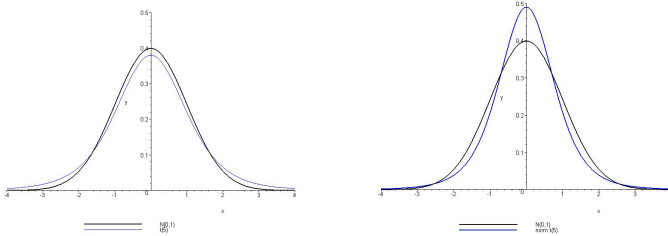
Značení: $t(\eta)$.

Hustota:

$$f_{t(\eta)}(x) = c(\eta) \left(1 + \frac{x^2}{\eta}\right)^{-\frac{1+\eta}{2}},$$
$$c(\eta) = \frac{\Gamma\left(\frac{1+\eta}{2}\right)}{\sqrt{\eta\pi}\Gamma\left(\frac{\eta}{2}\right)}.$$

Symetrie kolem nuly $\Rightarrow q_{t(\eta)}(1-\alpha) = -q_{t(\eta)}(\alpha)$.

Pro velký počet stupňů volnosti se nahrazuje normálním rozdělením.



Hustota normovaného normálního rozdělení a Studentova rozdělení s 5 stupni volnosti (původního a normovaného).

5.10.4 Odhad střední hodnoty při **neznámém** rozptylu II

V našem případě:

$$U = \frac{\sqrt{n}}{\sigma} (\bar{X} - \mu) \text{ má } N(0, 1),$$

$$V = \frac{(n-1) S_{\mathbf{X}}^2}{\sigma^2} \text{ má } \chi^2(n-1), \eta = n-1,$$

$$\frac{U}{\sqrt{\frac{V}{\eta}}} = \frac{\frac{\sqrt{n}}{\sigma} (\bar{X} - \mu)}{\sqrt{\frac{S_{\mathbf{X}}^2}{\sigma^2}}} = \frac{\sqrt{n}}{S_{\mathbf{X}}} (\bar{X} - \mu) \text{ má } t(n-1).$$

Z toho vyplývají intervalové odhady

$$\begin{aligned} & \left(-\infty, \bar{X} + \frac{S_{\mathbf{X}}}{\sqrt{n}} q_{t(n-1)}(1-\alpha) \right), \\ & \left(\bar{X} - \frac{S_{\mathbf{X}}}{\sqrt{n}} q_{t(n-1)}(1-\alpha), \infty \right), \\ & \left(\bar{X} - \frac{S_{\mathbf{X}}}{\sqrt{n}} q_{t(n-1)}\left(1 - \frac{\alpha}{2}\right), \bar{X} + \frac{S_{\mathbf{X}}}{\sqrt{n}} q_{t(n-1)}\left(1 - \frac{\alpha}{2}\right) \right). \end{aligned}$$

Při výpočtu nahradíme výběrový průměr \bar{X} jeho realizací \bar{x} a výběrovou směrodatnou odchylku $S_{\mathbf{X}}$ její realizací $s_{\mathbf{x}}$.

5.10.5 Odhad rozptylu

σ^2 odhadneme výběrovým rozptylem $S_{\mathbf{X}}^2$; $\frac{(n-1)S_{\mathbf{X}}^2}{\sigma^2}$ má rozdělení $\chi^2(n-1)$;

$$\begin{aligned} & P \left[\frac{(n-1) S_{\mathbf{X}}^2}{\sigma^2} \in (-\infty, q_{\chi^2(n-1)}(1-\alpha)) \right] \\ & = 1 - \alpha \\ & = P \left[\frac{(n-1) S_{\mathbf{X}}^2}{\sigma^2} \leq q_{\chi^2(n-1)}(1-\alpha) \right] \\ & = P \left[\frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}(1-\alpha)} \leq \sigma^2 \right] \\ & = P \left[\sigma^2 \in \left\langle \frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}(1-\alpha)}, \infty \right\rangle \right]. \end{aligned}$$

Dostali jsme **dolní** odhad.

Obdobně dostaneme i další intervalové odhady

$$\begin{aligned} & \left(-\infty, \frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}(\alpha)} \right), \\ & \left\langle \frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}(1-\alpha)}, \infty \right\rangle, \\ & \left\langle \frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}\left(1 - \frac{\alpha}{2}\right)}, \frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}\left(\frac{\alpha}{2}\right)} \right\rangle. \end{aligned}$$

Při výpočtu nahradíme výběrový rozptyl $S_{\mathbf{X}}^2$ jeho realizací $s_{\mathbf{x}}^2$.

5.10.6 Intervalové odhady spojitých rozdělání, která nejsou normální

převádíme obvykle na normální rozdělání nelineární transformací

$$h(t) = \Phi^{-1}(F_X(t))$$

($F_X(X)$ má rovnoměrné rozdělání na $\langle 0, 1 \rangle$).

Použijeme intervalový odhad pro normální rozdělání a transformujeme jej zpět podle vzorce

$$h^{-1}(u) = q_X^{-1}(\Phi(u)).$$

5.11 Obecné odhady parametrů

Rozdělání náhodné veličiny X závisí na vektoru parametrů $\boldsymbol{\vartheta} = (\vartheta_1, \dots, \vartheta_i) \in \Pi$, kde $\Pi \subseteq \mathbb{R}^i$ je **parametrický prostor**, tj. množina všech přípustných hodnot parametrů; pravděpodobnostní funkci značíme $p_X(t; \boldsymbol{\vartheta}) = p_X(t; \vartheta_1, \dots, \vartheta_i)$ atd.

Hledáme odhad $\hat{\boldsymbol{\Theta}} = (\hat{\Theta}_1, \dots, \hat{\Theta}_i)$, resp. realizaci odhadu $\hat{\boldsymbol{\vartheta}} = (\hat{\vartheta}_1, \dots, \hat{\vartheta}_i)$ pomocí realizace $\mathbf{x} = (x_1, \dots, x_n)$.

5.11.1 Metoda momentů

Pro $k = 1, 2, \dots$ je k -tý obecný moment funkcí $\boldsymbol{\vartheta}$,

$$EX^k(\boldsymbol{\vartheta}) = EX^k(\vartheta_1, \dots, \vartheta_i)$$

(závislost na parametrech lze stanovit dle pravděpodobnostního modelu).

Lze jej též odhadnout pomocí výběrového k -tého obecného momentu m_{X^k} .

Metoda momentů doporučuje realizaci odhadu $\hat{\boldsymbol{\vartheta}} = (\hat{\vartheta}_1, \dots, \hat{\vartheta}_i)$ takovou, že

$$EX^k(\hat{\vartheta}_1, \dots, \hat{\vartheta}_i) = m_{X^k} = \frac{1}{n} \sum_{j=1}^n x_j^k, \quad k = 1, 2, \dots$$

K jednoznačnému určení i proměnných obvykle potřebujeme (prvních) i rovnic pro $k = 1, 2, \dots, i$.

Použitelnost metody momentů

Možné problémy:

1. Řešení neexistuje \Rightarrow zkusme ubrat rovnice.
2. Je nekonečně mnoho řešení \Rightarrow zkusme přibrat další rovnice.
3. Je více než jedno řešení (např. soustavy kvadratických rovnic).
4. Je jediné řešení, ale je obtížné je nalézt.
5. Soustava je špatně podmíněná (typicky pro velký počet parametrů).
6. Našli jsme jediné řešení, které však **nesplňuje předpoklady**, $\hat{\boldsymbol{\vartheta}} \notin \Pi$ (např. parametry nemohou být libovolná čísla) \Rightarrow NELZE! **Vždy kontrolujte řešení!**

7. Všem rovnicím je přikládána stejná důležitost, což bývá nežádoucí (typicky pro velký počet parametrů).

8. Nelze použít pro nenumerická data (pokud je nelze smysluplně očíslovat).

Výhoda:

1. Lze použít pro diskrétní, spojitě i **smíšené** rozdělení beze změn.

5.11.2 Metoda maximální věrohodnosti (likelihood)

Pro diskrétní rozdělení Pravděpodobnost realizace,

$$\begin{aligned} p_{\mathbf{X}}(\mathbf{x}; \boldsymbol{\vartheta}) &= P[X_1 = x_1 \wedge \dots \wedge X_n = x_n; \boldsymbol{\vartheta}] \\ &= \prod_{j=1}^n P[X_j = x_j; \boldsymbol{\vartheta}] = \prod_{j=1}^n p_X(x_j; \boldsymbol{\vartheta}) = L(\boldsymbol{\vartheta}), \end{aligned}$$

je funkce $L: \Pi \rightarrow \langle 0, 1 \rangle$, $\Pi \subseteq \mathbb{R}^i$, parametrů $\boldsymbol{\vartheta} = (\vartheta_1, \dots, \vartheta_i)$, zvaná **věrohodnost realizace diskrétního rozdělení**.

Řešením jsou takové hodnoty $\hat{\boldsymbol{\vartheta}} = (\hat{\vartheta}_1, \dots, \hat{\vartheta}_i)$, které maximalizují věrohodnost. Maximalizujeme buď věrohodnost, nebo její logaritmus (*log-likelihood*),

$$\ell(\boldsymbol{\vartheta}) = \ln L(\boldsymbol{\vartheta}) = \sum_{j=1}^n \ln p_X(x_j; \boldsymbol{\vartheta}).$$

(Nutno vyloučit případ $p_X(x_j; \boldsymbol{\vartheta}) = 0$, který však nevede na maximum.)

Příklad: Empirické rozdělení je maximálně věrohodný odhad diskrétního rozdělení (pokud na rozdělení nejsou kladeny další podmínky).

Poznámka: Odhad na základě maxima věrohodnosti odpovídá Bayesovskému odhadu ve speciálním případě, kdy všechny hodnoty parametrů mají stejnou apriorní pravděpodobnost (resp. hustotu pravděpodobnosti). Používá se, pokud apriorní pravděpodobnosti parametrů neznáme.

Pro spojitě rozdělení

Každá realizace má nulovou pravděpodobnost, proto místo ní použijeme hustotu pravděpodobnosti, což ale vede na zcela **jiný pojem**

$$f_{\mathbf{X}}(\mathbf{x}; \boldsymbol{\vartheta}) = \prod_{j=1}^n f_X(x_j; \boldsymbol{\vartheta}) = \Lambda(\boldsymbol{\vartheta}).$$

Nicméně i tato funkce $\Lambda: \Pi \rightarrow \langle 0, \infty \rangle$, $\Pi \subseteq \mathbb{R}^i$, se nazývá **věrohodnost realizace spojitěho rozdělení**.

Pro korektní definici potřebujeme **spojitou** hustotu (alespoň na oboru hodnot, jichž náhodná veličina nabývá); taková hustota je nejvýše jedna.

$$\lambda(\boldsymbol{\vartheta}) = \ln \Lambda(\boldsymbol{\vartheta}) = \sum_{j=1}^n \ln f_X(x_j; \boldsymbol{\vartheta}).$$

(Nutno vyloučit případ $f_X(x_j; \boldsymbol{\vartheta}) = 0$, který však nevede na maximum.)

**Pro smíšené rozdělení
není věrohodnost definována!**

Použitelnost metody maximální věrohodnosti

Možné problémy:

1. Je více než jedno řešení. (Může se stát, že různé hodnoty parametrů popisují totéž rozdělení – vadí to?)
2. Řešení neexistuje (to se může stát jedině když věrohodnostní funkce je nespojitá nebo parametrický prostor neuzavřený).
3. Je jediné řešení, ale je obtížné je nalézt. (Lokální extrémy nemusí být globální.)
4. Soustava je špatně podmíněná.
5. Hodnoty věrohodnosti mohou být velmi malé.
6. **Nelze použít pro smíšené rozdělení!**

Výhody:

1. Hledání optima je o něco snazší než řešení soustavy rovnic.
2. Různým datům je dán společný (srovnatelný) význam.
3. Lze použít i na nenumerická data.

Příklad 2. Z realizace náhodného výběru $\mathbf{x} = (x_1, \dots, x_n)$ z normálního rozdělení $N(\mu, \sigma^2)$ odhadněte parametry μ a $r = \sigma^2$.

Řešení: Metoda momentů: Použijeme první dva obecné momenty,

$$EX = \mu, \quad EX^2 = (EX)^2 + DX = \mu^2 + \sigma^2 = \mu^2 + r.$$

Pro odhady $\hat{\mu}, \hat{r}$ máme soustavu rovnic

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^n x_j,$$
$$\hat{\mu}^2 + \hat{r} = \frac{1}{n} \sum_{j=1}^n x_j^2.$$

Řešení:

$$\hat{\mu} = \bar{\mathbf{x}},$$
$$\hat{r} = \frac{1}{n} \sum_{j=1}^n x_j^2 - \hat{\mu}^2 = \frac{1}{n} \sum_{j=1}^n x_j^2 - \bar{\mathbf{x}}^2,$$

což je alternativní (vychýlený konzistentní) odhad rozptylu

$$\begin{aligned}\widehat{DX} &= \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 = \frac{1}{n} \sum_{j=1}^n x_j^2 - \frac{2}{n} \bar{x} \sum_{j=1}^n x_j + \frac{1}{n} \sum_{j=1}^n \bar{x}^2 = \\ &= \frac{1}{n} \sum_{j=1}^n x_j^2 - 2 \bar{x}^2 + \bar{x}^2 = \frac{1}{n} \sum_{j=1}^n x_j^2 - \bar{x}^2 = \widehat{r}.\end{aligned}$$

Metoda maximální věrohodnosti:

$$\Lambda(\mu, r) = \prod_j f_{N(\mu, r)}(x_j) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi r}} \exp\left(-\frac{(x_j - \mu)^2}{2r}\right),$$

$$\lambda(\mu, r) = \ln \Lambda(\mu, r) = \frac{-1}{2r} \sum_{j=1}^n (x_j - \mu)^2 - \frac{n}{2} \ln r - \frac{n}{2} \ln 2\pi,$$

$$\frac{\partial}{\partial \mu} \lambda(\mu, r) = \frac{1}{r} \sum_{j=1}^n (x_j - \mu) = \frac{1}{r} \left(\sum_{j=1}^n x_j - n\mu \right) = \frac{n}{r} (\bar{x} - \mu),$$

$$\frac{\partial}{\partial r} \lambda(\mu, r) = \frac{1}{2r^2} \sum_{j=1}^n (x_j - \mu)^2 - \frac{n}{2r} = \frac{n}{2r^2} \left(\frac{1}{n} \sum_{j=1}^n (x_j - \mu)^2 - r \right).$$

Maximum opět nastává pro

$$\widehat{\mu} = \bar{x},$$

$$\widehat{r} = \frac{1}{n} \sum_{j=1}^n (x_j - \mu)^2 = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 = \widehat{DX}.$$

Odhad parametrů směsi normálních rozdělání

Úloha: Z realizace náhodného výběru (x_1, \dots, x_n) určete maximálně věrohodný odhad směsi normálních rozdělání se středními hodnotami μ_k , $k = 1, \dots, K$, stejným **známým** rozptylem σ^2 a koeficienty směsi (váhami) c_k , $k = 1, \dots, K$.

Pokus o řešení:

$$f_X(t) = \sum_k c_k f_{N(\mu_k, \sigma^2)}(t) = \sum_k c_k \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t - \mu_k)^2}{2\sigma^2}\right),$$

$$\begin{aligned}\Lambda(\mu, \mathbf{c}) &= \prod_j f_X(x_j) = \prod_j \sum_k c_k f_{N(\mu_k, \sigma^2)}(x_j) = \\ &= \prod_j \left(\frac{1}{\sqrt{2\pi}\sigma} \sum_k c_k \exp\left(-\frac{(x_j - \mu_k)^2}{2\sigma^2}\right) \right) = \\ &= \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \prod_j \sum_k c_k \exp\left(-\frac{(x_j - \mu_k)^2}{2\sigma^2}\right),\end{aligned}$$

$$\begin{aligned}\lambda(\mu, \mathbf{c}) &= \sum_j \ln \sum_k c_k f_{N(\mu_k, \sigma^2)}(x_j) = \\ &= -n \ln(\sqrt{2\pi}\sigma) + \sum_j \ln \sum_k c_k \exp\left(-\frac{(x_j - \mu_k)^2}{2\sigma^2}\right).\end{aligned}$$

Věrohodnost se těžko maximalizuje přímo, používá se iterační metoda:

EM algoritmus

EM (Expectation-Maximization) [Dempster, Laird, and Rubin 1977, M.I. Schlesinger 1968, US Army ~1950].

Stupeň příslušnosti x_j ke k -té složce směsi popíšeme koeficientem $\alpha_{j,k} \in (0, 1)$, přičemž

$$\sum_{k=1}^K \alpha_{j,k} = 1, \quad \sum_{j=1}^n \alpha_{j,k} > 0.$$

1. Zvolíme náhodně různé střední hodnoty složek směsi μ_k a nenulové koeficienty c_k , $k = 1, \dots, K$, splňující $\sum_k c_k = 1$.

E. Stanovíme stupně příslušnosti

$$\alpha_{j,k} := \frac{c_k f_{N(\mu_k, \sigma^2)}(x_j)}{\sum_{k'=1}^K c_{k'} f_{N(\mu_{k'}, \sigma^2)}(x_j)} = \frac{c_k \exp\left(\frac{-(x_j - \mu_k)^2}{2\sigma^2}\right)}{\sum_{k'=1}^K \left(c_{k'} \exp\left(\frac{-(x_j - \mu_{k'})^2}{2\sigma^2}\right)\right)}$$

(jmenovatel je normalizační faktor).

M. Aktualizujeme koeficienty složek směsi

$$c_k := \frac{\sum_{j=1}^n \alpha_{j,k}}{\sum_{k'=1}^K \sum_{j=1}^n \alpha_{j,k'}} = \frac{1}{n} \sum_{j=1}^n \alpha_{j,k}$$

a střední hodnoty složek jako těžiště hodnot realizace vážených stupni příslušnosti,

$$\mu_k := \frac{\sum_{j=1}^n \alpha_{j,k} x_j}{\sum_{j=1}^n \alpha_{j,k}} = \frac{\sum_{j=1}^n \alpha_{j,k} x_j}{n c_k}.$$

2. Opakujeme EM, dokud to přináší podstatnou změnu výsledků.

Podobně lze postupovat i pro neznámé rozptyly jednotlivých složek směsi.

Věta: V průběhu EM algoritmu **věrohodnost neklesá**.

Toto je jen velmi speciální ukázka EM algoritmu; lze jej snadno rozšířit na více dimenzí a jiné typy směsí.

Použití pro parametry směsí rozdělení je typické, ne však jediné možné.

Problém: Uvíznutí v lokálním extrému.

EM algoritmus rozšiřuje možnosti použití metody maximální věrohodnosti.

6 Testování hypotéz

6.1 Základní pojmy a principy testování hypotéz

(doporučená literatura: [Jaroš a kol.])

Máme posoudit hypotézu o hodnotě nějakého parametru rozdělení ϑ (pomocí **kritéria** čili **testovací statistiky** T , resp. její realizace t).

Předpoklad: Parametr ϑ nabývá pouze 2 hodnot, 0 pro „normální“ populaci, 1 pro „anomální“ prvky. O prvku máme rozhodnout, ke které skupině patří (tj. odhadnout ϑ). K tomu použijeme testovací statistiku T (resp. její realizaci t). Ta závisí na ϑ . Předpokládejme, že obě skupiny mají známá rozdělení statistiky T , která pro anomální skupinu nabývá „větších“ hodnot. (Některé hodnoty statistiky T se mohou vyskytnout v obou skupinách, takže klasifikace nemůže být bezchybná.) Zvolíme práh $\kappa \in \mathbb{R}$ a prvek klasifikujeme následovně:

pro $T \leq \kappa$ normální,

pro $T > \kappa$ anomální.

Příklad: Máme zastavit používání léku pro podezření z nežádoucích účinků?

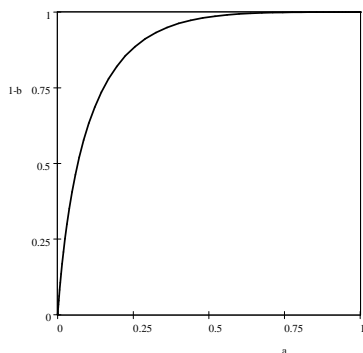
Nulová hypotéza H_0 : Výrobce je nevinný, riziko se nezvyšuje.

Alternativní hypotéza H_1 : Výrobce je vinný, riziko se zvyšuje.

Chyba 1. druhu (obviníme nevinného): Zamítneme nulovou hypotézu, která platí. Normální je klasifikován jako anomální s pravděpodobností $\alpha(\kappa)$ (nerostoucí funkce κ).

Chyba 2. druhu (osvobodíme vinného): Nezamítneme nulovou hypotézu, která neplatí. Anomální je klasifikován jako normální s pravděpodobností $\beta(\kappa)$ (neklesající funkce κ).

ROC křivka (angl. **ROC curve, receiver operating characteristic**) vyjadřuje závislost pravděpodobnosti chyby prvního druhu α (vodorovně) a síly testu $1 - \beta$ (svisle), parametrem křivky je kritická hodnota κ . Volbou kritické hodnoty se chceme co nejvíce přiblížit bodu $(0, 1)$, tj. bezchybné klasifikaci. Nicméně vybereme bod, v němž se pravděpodobnost chyby prvního druhu rovná zvolenému číslu α (tj. s danou vodorovnou souřadnicí).



Obrázek 1: Typický průběh **ROC** křivky

Možná kritéria pro volbu prahu κ :

- $\alpha(\kappa) = \beta(\kappa)$,
- $\min_{\kappa}(\alpha(\kappa) + \beta(\kappa))$,
- $\min_{\kappa} e(\alpha(\kappa), \beta(\kappa))$, např. $\min_{\kappa}(a\alpha(\kappa) + b\beta(\kappa))$, tj. minimalizace **výplatní funkce**,
- $\alpha(\kappa) =$ předem zvolená malá hodnota.

Většinou se používá poslední možnost, a to z důvodů

- technických (snazší úloha),
- nepotřebujeme znát rozdělení anomální skupiny,
- obvykle máme více než dvě možné hodnoty parametru, což situaci komplikuje.

Volbou přísnosti kritéria snižujeme riziko jedné chyby na úkor zvýšení rizika druhé chyby.

Dohodnuté východisko: **Kritickou hodnotu** testu κ stanovíme tak, aby chyba 1. druhu nastávala s danou pravděpodobností α zvanou **hladina významnosti** (nebo s menší pravděpodobností, nelze-li dosáhnout rovnost).

Podle tradice v oboru se nejčastěji užívají hodnoty 1% nebo 5% (vždy $\alpha \ll \frac{1}{2}$).

Hodnoty kritéria, která přesahují kritickou hodnotu (odpovídají výsledkům málo pravděpodobným při platnosti nulové hypotézy) považujeme za **statisticky významné** a v tom případě **nulovou hypotézu zamítáme**.

V opačném případě **nulovou hypotézu nezamítáme**, ale **ani nepotvrzujeme**, neboť tím bychom se mohli dopustit chyby 2. druhu s blíže neurčenou pravděpodobností β .

Síla testu $1 - \beta$.

Rozlišuje se

- **jednoduchá hypotéza**: nulové hypotéze odpovídá jediná hodnota parametru,
- **složená hypotéza**: nulové hypotéze odpovídá více hodnot parametru,

a dále

- **jednoduchá alternativa**: alternativní hypotéze odpovídá jediná hodnota parametru,
- **složená alternativa**: alternativní hypotéze odpovídá více hodnot parametru.

Často se formuluje nulová a alternativní hypotéza tak, že nejsou navzájem svými negacemi a nepokrývají prostor všech možných hodnot parametru. Vzniká tím jen chaos (viz většina ostatní literatury). Snadno se mu vyhneme, když budeme formulovat nulovou hypotézu jako negaci alternativní hypotézy.

Je-li např. $H_1 : \vartheta > c$, pak nevolíme $H_0 : \vartheta = c$, ale $H_0 : \vartheta \leq c$. (Největší riziko chyby 1. druhu obvykle odpovídá případu $\vartheta = c$, takže postup je stejný.)

U složené hypotézy požadujeme, aby pravděpodobnost chyby 1. druhu byla nejvýše α pro všechny hodnoty parametru vyhovující nulové hypotéze.

(Statistická významnost neznamená významnost praktickou.)

Řešení: Nulovou hypotézu zamítneme, právě když hodnota kritéria získaná z realizace nezapadne do intervalu spolehlivosti pro koeficient spolehlivosti $1 - \alpha$, tj. kritická hodnota je mezi intervalového odhadu.

Obrácený problém: Při jaké mezní hladině významnosti by pozorovaná hodnota byla kritická; tomu říkáme **dosažená významnost**; stačí ji porovnat s předem zvolenou hladinou významnosti testu. (**Čím nižší číslo, tím významnější výsledek.**) Programy obvykle dávají za výsledek dosaženou významnost (obvykle se značí P a říká se jí pouze *significance*). Výhody: hladinu významnosti není třeba předem zadat, a navíc se dovíme, jak daleko od ní jsme byli.

Typický tvar testu: Testovací statistiku T , která roste s parametrem ϑ a má známé rozdělení, (přesněji její realizaci t) porovnáváme s kvantily příslušného rozdělení a zamítneme při extrémních hodnotách (nepravděpodobných při platnosti nulové hypotézy):

H_0	H_1	zamítáme pro	dosažená významnost
$\vartheta \leq c$	$\vartheta > c$	$t > q_T(1 - \alpha)$	$1 - F_T(t)$
$\vartheta \geq c$	$\vartheta < c$	$t < q_T(\alpha)$	$F_T(t)$
$\vartheta = c$	$\vartheta \neq c$	$t > q_T(1 - \frac{\alpha}{2})$ nebo $t < q_T(\frac{\alpha}{2})$	$2 \min(F_T(t), 1 - F_T(t))$

V literatuře se setkáme i s následujícími případy hypotéz, které se však řeší stejně jako první dva výše uvedené:

H_0	H_1
$\vartheta = c$	$\vartheta > c$
$\vartheta = c$	$\vartheta < c$

6.2 Testy střední hodnoty normálního rozdělení

6.2.1 Při známém rozptylu σ^2

$$t = \frac{\bar{x} - c}{\sigma} \sqrt{n}$$

porovnáváme s kvantily **normovaného normálního rozdělení**:

H_0	zamítáme pro	dosažená významnost
$\mu \leq c$	$t > \Phi^{-1}(1 - \alpha)$	$1 - \Phi(t)$
$\mu \geq c$	$t < -\Phi^{-1}(1 - \alpha) = \Phi^{-1}(\alpha)$	$\Phi(t)$
$\mu = c$	$ t > \Phi^{-1}(1 - \frac{\alpha}{2})$	$2(1 - \Phi(t))$

6.2.2 Při neznámém rozptylu

$$t = \frac{\bar{x} - c}{s_x} \sqrt{n}$$

porovnáváme s kvantily **Studentova rozdělení** s $n - 1$ stupni volnosti:

H_0	zamítáme pro	dosažená významnost
$\mu \leq c$	$t > q_{t(n-1)}(1 - \alpha)$	$1 - F_{t(n-1)}(t)$
$\mu \geq c$	$t < -q_{t(n-1)}(1 - \alpha)$	$F_{t(n-1)}(t)$
$\mu = c$	$ t > q_{t(n-1)}(1 - \frac{\alpha}{2})$	$2(1 - F_{t(n-1)}(t))$

6.3 Testy rozptylu normálního rozdělení

$$t = \frac{(n-1)s_x^2}{c}$$

porovnááme s kvantily χ^2 -**rozdělení** s $n-1$ stupni volnosti:

H_0	zamítáme pro	dosažená významnost
$\sigma^2 \leq c$	$t > q_{\chi^2(n-1)}(1-\alpha)$	$1 - F_{\chi^2(n-1)}(t)$
$\sigma^2 \geq c$	$t < q_{\chi^2(n-1)}(\alpha)$	$F_{\chi^2(n-1)}(t)$
$\sigma^2 = c$	$t < q_{\chi^2(n-1)}(\frac{\alpha}{2})$ nebo $t > q_{\chi^2(n-1)}(1 - \frac{\alpha}{2})$	$2 \min(F_{\chi^2(n-1)}(t), 1 - F_{\chi^2(n-1)}(t))$

6.4 Porovnání dvou normálních rozdělení

Předpoklad: **Nezávislé** výběry

(X_1, \dots, X_m) z rozdělení $N(EX, DX)$,

(Y_1, \dots, Y_n) z rozdělení $N(EY, DY)$.

6.4.1 Testy rozptylu dvou normálních rozdělení [Fisher]

Je-li $DX = DY$, pak $S_X^2 \doteq S_Y^2$. Testovací statistikou je

$$T = \frac{S_X^2}{S_Y^2}.$$

F-rozdělení (Fisherovo-Snedecorovo rozdělení) s ξ a η stupni volnosti je rozdělení náhodné veličiny

$$F = \frac{\frac{U}{\xi}}{\frac{V}{\eta}},$$

kde U, V jsou **nezávislé** náhodné veličiny s rozdělením $\chi^2(\xi)$, resp. $\chi^2(\eta)$.

Značení: $F(\xi, \eta)$

Hustota pro $x > 0$:

$$f_{F(\xi, \eta)}(x) = c(\xi, \eta) x^{\frac{\xi}{2}-1} \left(1 + \frac{\xi}{\eta} x\right)^{-\frac{\xi+\eta}{2}},$$

$$c(\xi, \eta) = \frac{\Gamma\left(\frac{\xi+\eta}{2}\right)}{\Gamma\left(\frac{\xi}{2}\right) \Gamma\left(\frac{\eta}{2}\right)} \left(\frac{\xi}{\eta}\right)^{\frac{\xi}{2}}$$

Je-li $DX = DY = \sigma^2$, pak dosadíme

$$U := \frac{(m-1)S_X^2}{\sigma^2} \text{ má } \chi^2(m-1),$$

$$V := \frac{(n-1)S_Y^2}{\sigma^2} \text{ má } \chi^2(n-1),$$

$$\xi := m-1, \eta := n-1,$$

$$F = \frac{\frac{U}{\xi}}{\frac{V}{\eta}} = \frac{\frac{(m-1)S_X^2}{(m-1)\sigma^2}}{\frac{(n-1)S_Y^2}{(n-1)\sigma^2}} = \frac{S_X^2}{S_Y^2} = T.$$

Testujeme realizaci

$$t = \frac{s_x^2}{s_y^2}$$

na rozdělení $F(m-1, n-1)$:

H_0	zamítáme pro	dosažená významnost
$DX \leq DY$	$t > q_{F(m-1, n-1)}(1 - \alpha)$	$1 - F_{F(m-1, n-1)}(t)$
$DX \geq DY$	$t < q_{F(m-1, n-1)}(\alpha)$	$F_{F(m-1, n-1)}(t)$
$DX = DY$	$t < q_{F(m-1, n-1)}(\frac{\alpha}{2})$ nebo $t > q_{F(m-1, n-1)}(1 - \frac{\alpha}{2})$	$2 \min(F_{F(m-1, n-1)}(t), 1 - F_{F(m-1, n-1)}(t))$

Pro každou hladinu významnosti potřebujeme dvoudimenzionální tabulku kvantilů indexovanou ξ, η ; obvykle je tabelována jen polovina, druhou je třeba dopočítat podle vzorce

$$q_{F(\xi, \eta)}(\beta) = \frac{1}{q_{F(\eta, \xi)}(1 - \beta)}.$$

(Pozor na opačné pořadí indexů!)

Lépe je uvažovat $\frac{s_y^2}{s_x^2}$ místo $\frac{s_x^2}{s_y^2}$, takže rozlišíme 2 případy:

1. Pro $s_x^2 \geq s_y^2$ testujeme

$$t = \frac{s_x^2}{s_y^2} \geq 1$$

na rozdělení $F(m-1, n-1)$:

H_0	zamítáme pro	dosažená významnost
$DX \leq DY$	$t > q_{F(m-1, n-1)}(1 - \alpha)$	$1 - F_{F(m-1, n-1)}(t)$
$DX \geq DY$	nezamítáme	žádná
$DX = DY$	$t > q_{F(m-1, n-1)}(1 - \frac{\alpha}{2})$	$2 (1 - F_{F(m-1, n-1)}(t))$

1. Pro $s_x^2 \leq s_y^2$ testujeme

$$t = \frac{s_y^2}{s_x^2} \geq 1$$

na rozdělení $F(n-1, m-1)$ (pozor na pořadí počtů stupňů volnosti!):

H_0	zamítáme pro	dosažená významnost
$DX \leq DY$	nezamítáme	žádná
$DX \geq DY$	$t > q_{F(n-1, m-1)}(1 - \alpha)$	$1 - F_{F(n-1, m-1)}(t)$
$DX = DY$	$t > q_{F(n-1, m-1)}(1 - \frac{\alpha}{2})$	$2 (1 - F_{F(n-1, m-1)}(t))$

6.4.2 Testy středních hodnot dvou normálních rozdělení se známým rozptylem σ^2

$$\bar{X}_m \text{ má } N\left(EX, \frac{\sigma^2}{m}\right),$$

$$\bar{Y}_n \text{ má } N\left(EY, \frac{\sigma^2}{n}\right),$$

$$\bar{X}_m - \bar{Y}_n \text{ má } N\left(EX - EY, \sigma^2 \left(\frac{1}{m} + \frac{1}{n}\right)\right).$$

Za předpokladu $EX = EY$:

$$T := \frac{\bar{X}_m - \bar{Y}_n}{\sigma \sqrt{\frac{1}{m} + \frac{1}{n}}} \text{ má } N(0, 1).$$

Testujeme realizaci t na $N(0, 1)$ (viz kapitola 6.2.1).

6.4.3 Testy středních hodnot dvou normálních rozdělení se (stejným) **neznámým** rozptylem

Předpoklad: $DX = DY = \sigma^2$

Nejprve ověříme tento předpoklad (viz kapitola 6.4.1).

(Ve skutečnosti nemůžeme předpoklad ověřit, jedinečně vyvrátit; pokusíme se o to, a pokud se to nepodaří, pokračujeme. Bez tohoto předpokladu by byl další postup složitější, viz např. [Mood a kol.]).

Máme dva odhady $S_{\mathbf{X}}^2, S_{\mathbf{Y}}^2$ stejné hodnoty σ^2 ; použijeme jejich průměr vážený rozsahy výběrů (-1 kvůli výpočtu výběrového průměru):

$$\begin{aligned} \frac{(m-1)S_{\mathbf{X}}^2}{\sigma^2} & \text{ má } \chi^2(m-1), \\ \frac{(n-1)S_{\mathbf{Y}}^2}{\sigma^2} & \text{ má } \chi^2(n-1), \\ \frac{(m-1)S_{\mathbf{X}}^2 + (n-1)S_{\mathbf{Y}}^2}{\sigma^2} & \text{ má } \chi^2(m+n-2) \end{aligned}$$

se střední hodnotou $m+n-2$,

$$\frac{(m-1)S_{\mathbf{X}}^2 + (n-1)S_{\mathbf{Y}}^2}{(m+n-2)\sigma^2} = \frac{S^2}{\sigma^2}$$

má střední hodnotu 1 a

$$S^2 := \frac{(m-1)S_{\mathbf{X}}^2 + (n-1)S_{\mathbf{Y}}^2}{m+n-2}$$

je nestranný odhad σ^2 ,

$$S := \sqrt{\frac{(m-1)S_{\mathbf{X}}^2 + (n-1)S_{\mathbf{Y}}^2}{m+n-2}}.$$

$$\bar{X}_m \text{ má } N\left(EX, \frac{\sigma^2}{m}\right),$$

$$\bar{Y}_n \text{ má } N\left(EY, \frac{\sigma^2}{n}\right),$$

$$\bar{X}_m - \bar{Y}_n \text{ má } N\left(EX - EY, \sigma^2\left(\frac{1}{m} + \frac{1}{n}\right)\right).$$

Za předpokladu $EX = EY$:

$$\frac{\bar{X}_m - \bar{Y}_n}{\sigma \sqrt{\frac{1}{m} + \frac{1}{n}}} \text{ má } N(0, 1),$$

$$\frac{(m+n-2)S^2}{\sigma^2} = \frac{(m-1)S_X^2 + (n-1)S_Y^2}{\sigma^2} \text{ má } \chi^2(m+n-2),$$

$$T := \frac{\bar{X}_m - \bar{Y}_n}{S \sqrt{\frac{1}{m} + \frac{1}{n}}} = \frac{\frac{\bar{X}_m - \bar{Y}_n}{\sigma \sqrt{\frac{1}{m} + \frac{1}{n}}}}{\sqrt{\frac{S^2}{\sigma^2}}} \text{ má } t(m+n-2).$$

Testujeme realizaci t na rozdělení $t(m+n-2)$ (viz kapitola 6.2.2).

6.5 Testy středních hodnot dvou normálních rozdělení - párový pokus

(dle [SH10])

Příklad: Máme porovnat průměrnou teplotu na dvou místech.

Standardní test středních hodnot dvou normálních rozdělení je slabý kvůli velkému rozptylu, který však má společnou příčinu a projevuje se proto synchronně v obou výběrech; proto výběry **nejsou navzájem nezávislé**. Měříme vždy obě veličiny současně.

Předpoklad: Náhodné veličiny X_j, Y_j ($j = 1, \dots, n$) mají normální rozdělení $N(\mu_j, \sigma^2)$ se stálým rozptylem σ^2 a proměnnými středními hodnotami $\mu_j = EX_j = EY_j$.

Můžeme použít náhodné veličiny $U_j := X_j - \mu_j, V_j := Y_j - \mu_j$ ($j = 1, \dots, n$), které **jsou nezávislé** a mají rozdělení $N(0, \sigma^2)$.

Náhodné veličiny $\Delta_j := X_j - Y_j = U_j - V_j$ ($j = 1, \dots, n$) jsou nezávislé a mají rozdělení $N(0, 2\sigma^2)$.

Výběrový průměr $\bar{\Delta}$ má $N\left(0, \frac{2\sigma^2}{n}\right)$.

6.5.1 Pro známý rozptyl σ^2

Neznámé parametry sdruženého rozdělení jsou μ_1, \dots, μ_n , ale nepotřebujeme je.

Dle kapitoly 6.2.1 (pro $c = 0$) testujeme

$$T := \frac{\bar{\Delta}}{\sigma} \sqrt{\frac{n}{2}} = \frac{\bar{X} - \bar{Y}}{\sigma} \sqrt{\frac{n}{2}}$$

na $N(0, 1)$.

6.5.2 Pro neznámý rozptyl

Neznámé parametry sdruženého rozdělení jsou $\Theta = (\sigma^2, \mu_1, \dots, \mu_n)$, potřebujeme z nich pouze $\sigma^2 = DX$.

Můžeme pracovat přímo s výběrem $(\Delta_1, \dots, \Delta_n)$ z normálního rozdělení.

Dle kapitoly 6.2.2 (pro $c = 0$) testujeme

$$T := \frac{\bar{\Delta}}{S_{\Delta}} \sqrt{n}$$

na $t(n-1)$.

Cvičení: Maximálně věrohodný odhad parametrů:

$$\begin{aligned}\ell(\Theta) &= \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma^2}\right) \cdot \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_j - \mu_j)^2}{2\sigma^2}\right), \\ L(\Theta) &= -\sum_{j=1}^n \frac{(x_j - \mu_j)^2}{2\sigma^2} - \sum_{j=1}^n \frac{(y_j - \mu_j)^2}{2\sigma^2} - 2n \ln \sigma - 2n \ln \sqrt{2\pi}, \\ 0 = \frac{\partial L(\hat{\Theta})}{\partial \hat{\mu}_j} &= \frac{\partial}{\partial \hat{\mu}_j} \left(-\frac{(x_j - \hat{\mu}_j)^2}{2\hat{\sigma}^2} - \frac{(y_j - \hat{\mu}_j)^2}{2\hat{\sigma}^2} \right) \\ &= \frac{1}{\hat{\sigma}^2} ((x_j - \hat{\mu}_j) + (y_j - \hat{\mu}_j)) = \frac{1}{\hat{\sigma}^2} (x_j + y_j - 2\hat{\mu}_j), \\ \hat{\mu}_j &= \frac{x_j + y_j}{2}, \quad j = 1, \dots, n.\end{aligned}$$

Odhady $\hat{\mu}_j$, ($j = 1, \dots, n$) **nejso konzistentní**.

Po jejich dosažení:

$$\begin{aligned}L(\hat{\Theta}) &= -\sum_{j=1}^n \frac{(x_j - y_j)^2}{4\hat{\sigma}^2} - n \ln \hat{\sigma}^2 - 2n \ln \sqrt{2\pi}, \\ 0 = \frac{\partial L(\hat{\Theta})}{\partial (\hat{\sigma}^2)} &= \sum_{j=1}^n \frac{(x_j - y_j)^2}{4(\hat{\sigma}^2)^2} - \frac{2n}{\hat{\sigma}^2}, \\ \hat{\sigma}^2 &= \frac{1}{2n} \sum_{j=1}^n (x_j - y_j)^2 = \frac{1}{2n} \sum_{j=1}^n \delta_j^2,\end{aligned}$$

kde δ_j je realizace Δ_j . Odhad $\hat{\sigma}^2$ **je konzistentní**.

6.6 χ^2 -test dobré shody

Slouží k testování hypotézy, že náhodná veličina má předpokládané rozdělení. Protože umíme hypotézy jen zamítat, nikdy nepotvrdíme, že takové rozdělení opravdu má.

Testujeme **diskrétní rozdělení** (mohlo vzniknout diskretizací spojitého).

H_0 : Náhodná veličina má diskrétní rozdělení do k tříd s nenulovými pravděpodobnostmi p_1, \dots, p_k .

Testujeme pomocí realizace náhodného výběru rozsahu n . Není důležité pořadí výsledků, pouze jejich **četnosti** n_i , resp. **relativní četnosti** $\frac{n_i}{n}$ ($i = 1, \dots, k$). Porovnáváme četnost n_i s **teoretickou četností** $n p_i$. Testovací statistikou je

$$T := \sum_{i=1}^k \frac{(n_i - n p_i)^2}{n p_i}.$$

Její rozdělení se pro $n \rightarrow \infty$ blíží $\chi^2(k-1)$.

Dosažená významnost: $1 - F_{\chi^2(k-1)}(t)$. Nulovou hypotézu zamítáme pro $t > q_{\chi^2(k-1)}(1 - \alpha)$, tj. $1 - F_{\chi^2(k-1)}(t) < \alpha$.

Cvičení 1. Tabulka udává rozdělení (podmíněné) pravděpodobnosti, že volič strany zastoupené v parlamentu volil danou stranu. Posuďte na 5% hladině významnosti hypotézu, že stejné rozdělení mají i poslanci.

relativní preference	0.376	0.344	0.136	0.077	0.067
počet poslanců	81	74	26	13	6

Řešení. Doplníme tabulku (poslední sloupec uvádí celkový údaj):

relativní preference	0.376	0.344	0.136	0.077	0.067	1
počet poslanců	81	74	26	13	6	200
teor. četnost	75.2	68.8	27.2	15.4	13.4	200
příspěvek k χ^2	0.447	0.393	0.052	0.374	4.086	5.353

Hodnotu kritéria 5.353 porovnáme s kvantilem $q_{\chi^2(4)}(0.95) \doteq 9.4877$ a hypotézu nezamítáme (poněkud překvapivý závěr vzhledem k tomu, že poslední dvě strany mají téměř stejnou podporu u voličů, ale poslední má více než $2\times$ méně poslanců).

6.6.1 Modifikace

Problém: Testujeme na rozdělení, kterému se skutečné jen limitně blíží. Tím se dopouštíme blíže neurčené dodatečné chyby. Teoretické četnosti tříd nesmí být příliš malé (řekněme aspoň 5), aby náš předpoklad byl oprávněný.

Modifikace: Vychází-li teoretická četnost některých tříd příliš malá, sloučíme je s jinými třídami (pokud možno „blízkými“).

Problém: Zkoumané rozdělení může záviset na neznámých parametrech.

Modifikace 1: Parametry odhadneme na základě **jiného** náhodného výběru.

Modifikace 2: Parametry odhadneme na základě **stejného** náhodného výběru, který používáme k testu dobré shody. Tím jsme však snížili počet stupňů volnosti, takže musíme testovat na rozdělení $\chi^2(k-1-q)$, kde q je počet odhadnutých parametrů.

Problém: Chceme testovat shodu se **spojitým** nebo **smíšeným** rozdělením.

Modifikace: Rozdělení napřed diskretizujeme, tj. všechny možné výsledky rozdělíme do k disjunktních tříd. Prvky v jedné třídě si mají být „blízké“, jinak snižujeme sílu testu. Všechny teoretické četnosti musí být dostatečně velké a nejlépe zhruba stejné.

Poznámka: Zásadně musíme pracovat s jednotkami (objekty), z nichž každá zvlášť (a nezávisle) je zařazena do nějaké třídy. Nelze počítat s tisíci, procenty, spojitým množstvím atd.

6.6.2 χ^2 -test dobré shody dvou rozdělení

(dle [Mood a kol.])

H_0 : Dvě diskrétní náhodné veličiny mají stejné diskrétní rozdělení.

Rozsahy výběrů jsou m, n a četnosti výsledků m_i, n_i ($i = 1, \dots, k$). Předpokládáme rozdělení

s neznámými teoretickými pravděpodobnostmi p_i ($i = 1, \dots, k$).

$$\sum_{i=1}^k \frac{(m_i - m p_i)^2}{m p_i} \text{ se blíží } \chi^2(k-1),$$

$$\sum_{i=1}^k \frac{(n_i - n p_i)^2}{n p_i} \text{ se blíží } \chi^2(k-1),$$

$$T = \sum_{i=1}^k \frac{(m_i - m p_i)^2}{m p_i} + \sum_{i=1}^k \frac{(n_i - n p_i)^2}{n p_i} \text{ se blíží } \chi^2(2(k-1)).$$

Neznámé parametry p_i odhadneme pomocí maxima věrohodnosti,

$$p_i = \frac{m_i + n_i}{m + n},$$

z nich je jen $k-1$ nezávislých (neboť $\sum_{i=1}^k p_i = 1$), takže výsledný počet stupňů volnosti je $2(k-1) - (k-1) = k-1$ a testujeme T na $\chi^2(k-1)$. Nulovou hypotézu zamítáme pro $t > q_{\chi^2(k-1)}(1-\alpha)$, tj. $1 - F_{\chi^2(k-1)}(t) < \alpha$. Praktičtější (ekvivalentní) vzorec:

$$T = \left(\frac{1}{m} + \frac{1}{n} \right) \sum_{i=1}^k \frac{(m_i - m p_i)^2}{p_i}.$$

6.6.3 χ^2 -test nezávislosti dvou rozdělení

(dle [Líkaš, Machek])

H_0 : Dvě náhodné veličiny (jejichž rozdělení neznáme) jsou nezávislé.

X nabývá k hodnot s pravděpodobnostmi p_1, \dots, p_k ,

Y nabývá m hodnot s pravděpodobnostmi q_1, \dots, q_m .

Realizace dvojrozměrného náhodného výběru $((x_1, y_1), \dots, (x_n, y_n))$ obsahuje dvojice realizací náhodných veličin X, Y ; z výsledků nás zajímají opět pouze četnosti n_{ij} ($i = 1, \dots, k$; $j = 1, \dots, m$). Ty bývají uspořádány do tzv. **kontingenční tabulky**. Počet tříd je km .

Za předpokladu nezávislosti jsou pravděpodobnosti výsledků $p_i q_j$ ($i = 1, \dots, k$; $j = 1, \dots, m$),

$$T := \sum_{i=1}^k \sum_{j=1}^m \frac{(n_{ij} - n p_i q_j)^2}{n p_i q_j} \text{ se blíží } \chi^2(km-1).$$

Neznámé parametry p_i, q_j odhadneme pomocí maxima věrohodnosti,

$$p_i = \frac{\sum_{j=1}^m n_{ij}}{n}, \quad q_j = \frac{\sum_{i=1}^k n_{ij}}{n},$$

z nich je jen $(k-1) + (m-1)$ nezávislých (neboť $\sum_{i=1}^k p_i = 1$, $\sum_{j=1}^m q_j = 1$), takže výsledný počet stupňů volnosti je $km - 1 - (k-1) - (m-1) = (k-1)(m-1)$ a testujeme T na $\chi^2((k-1)(m-1))$. Nulovou hypotézu zamítáme pro $t > q_{\chi^2((k-1)(m-1))}(1-\alpha)$, tj. $1 - F_{\chi^2((k-1)(m-1))}(t) < \alpha$.

6.7 Korelace, její odhad a testování

(dle [Líkaš, Machek])

Korelace $\rho(X, Y)$ náhodných veličin X, Y (s nenulovým rozptylem) je střední hodnota součinu odpovídajících normovaných veličin $\frac{X - EX}{\sigma_X} \cdot \frac{Y - EY}{\sigma_Y}$,

$$\rho(X, Y) = \frac{E((X - EX)(Y - EY))}{\sigma_X \sigma_Y} \in \langle -1, 1 \rangle.$$

Je nulová pro nezávislé náhodné veličiny, ale i pro některé jiné, tzv. **nekorelované**.

Extrémní hodnoty ± 1 odpovídají lineární závislosti mezi X, Y .

Na základě dvojrozměrného náhodného výběru $((X_1, Y_1), \dots, (X_n, Y_n))$ můžeme korelaci odhadnout pomocí **výběrového koeficientu korelace**

$$R_{\mathbf{X}, \mathbf{Y}} = \frac{\sum_{j=1}^n (X_j - \bar{\mathbf{X}})(Y_j - \bar{\mathbf{Y}})}{\sqrt{\left(\sum_{j=1}^n (X_j - \bar{\mathbf{X}})^2\right) \left(\sum_{j=1}^n (Y_j - \bar{\mathbf{Y}})^2\right)}}.$$

Jeho realizace

$$r_{\mathbf{x}, \mathbf{y}} = \frac{\sum_{j=1}^n (x_j - \bar{\mathbf{x}})(y_j - \bar{\mathbf{y}})}{\sqrt{\left(\sum_{j=1}^n (x_j - \bar{\mathbf{x}})^2\right) \left(\sum_{j=1}^n (y_j - \bar{\mathbf{y}})^2\right)}} \in \langle -1, 1 \rangle,$$

neboť je to kosinus úhlu vektorů

$$(x_1 - \bar{\mathbf{x}}, \dots, x_n - \bar{\mathbf{x}}), (y_1 - \bar{\mathbf{y}}, \dots, y_n - \bar{\mathbf{y}}) \in \mathbb{R}^n$$

neboli korelace empirického rozdělení, $r_{\mathbf{x}, \mathbf{y}} = \rho(\text{Emp}(\mathbf{x}, \mathbf{y}))$.

Pro výpočet se používá jednorůchodový vzorec:

$$R_{\mathbf{X}, \mathbf{Y}} = \frac{n \sum_{j=1}^n x_j y_j - \left(\sum_{j=1}^n x_j\right) \left(\sum_{j=1}^n y_j\right)}{\sqrt{\left(n \sum_{j=1}^n x_j^2 - \left(\sum_{j=1}^n x_j\right)^2\right) \left(n \sum_{j=1}^n y_j^2 - \left(\sum_{j=1}^n y_j\right)^2\right)}}.$$

6.7.1 Test nekorelovanosti dvou **normálních** rozdělení

Předpoklad: Dvojrozměrná náhodná veličina (X, Y) má (dvojrozměrné) normální rozdělení, $n \geq 3$.

H_0 : $\rho(X, Y) = 0$ (X, Y jsou nekorelované).

Testovací statistikou je

$$T = \frac{R_{\mathbf{X}, \mathbf{Y}} \sqrt{n-2}}{\sqrt{1 - R_{\mathbf{X}, \mathbf{Y}}^2}},$$

za předpokladu nekorelovanosti má rozdělení $t(n-2)$, dále postupujeme dle kapitoly 6.2.2.

6.8 Neparametrické testy

Jsou použitelné bez ohledu na typ rozdělení, jsou však slabší.

6.8.1 Znaménkový test

Rozlišujeme pouze znaménko odchylky od zvolené hodnoty c . Tím ztrácíme kvantativní informaci a tedy i možnost testovat např. střední hodnotu. Místo ní testujeme medián $q_X(\frac{1}{2})$.

$$H_0 : q_X(\frac{1}{2}) = c$$

Při platnosti nulové hypotézy by kladné i záporné odchylky měly být stejně pravděpodobné. Nulové odchylky z výběru předem vyloučíme. Testovací statistikou T je počet kladných odchylek, který testujeme na binomické rozdělení $\text{Bin}(n, \frac{1}{2})$. Nulovou hypotézu zamítáme pro

$$t < q_{\text{Bin}(n, \frac{1}{2})} \left(\frac{\alpha}{2} \right) \text{ nebo } t > q_{\text{Bin}(n, \frac{1}{2})} \left(1 - \frac{\alpha}{2} \right).$$

(Podobně pro jednostranné testy.) Výpočet kvantilů je pracný, ale kritické hodnoty jsou tabulovány (v závislosti na n a hladině významnosti).

Dosažená významnost se počítá o trochu snáze.

Pro velká n používáme centrální limitní větu a testujeme

$$T_0 := \frac{2T - n}{\sqrt{n}}$$

na $N(0, 1)$.

Lze použít i k porovnání dvou mediánů u párového pokusu.

Příklad použití: Odhad smrtelné dávky látky.

Na rozdíl od střední hodnoty medián vždy existuje (je však problém, jak ho definovat, aby byl jednoznačný).

Jeho výpočetní složitost je větší, řádu $n \ln n$.

6.8.2 Wilcoxonův test (jednovýběrový)

$H_0 : X$ má rozdělení symetrické kolem hodnoty c

(V tom případě je c mediánem i střední hodnotou.)

Z realizace (x_1, \dots, x_n) vypočteme posloupnost (z_1, \dots, z_n) , kde $z_j = x_j - c$. Seřadíme ji vzestupně podle absolutních hodnot $|z_j| = |x_j - c|$, čímž j -tému prvku přiřadíme pořadí r_j . Je-li více stejných rozdílů, přiřadíme jim stejné pořadí rovné aritmetickému průměru. Testovací statistikou je

$$T_1 := \sum_{j: z_j > 0} r_j$$

nebo

$$T_2 := \min \left(\sum_{j: z_j > 0} r_j, \sum_{j: z_j < 0} r_j \right),$$

porovnáme s tabulkou kritických hodnot pro tento test.

7 Co zde nebylo

- 7.1 Více o zobrazení náhodné veličiny funkcí a o součtu náhodných veličin
- 7.2 Diskretizace
- 7.3 Směs pravděpodobností
- 7.4 Charakteristická funkce náhodné veličiny
- 7.5 Důkaz centrální limitní věty

Literatura

- [Navara: PMS] Navara, M.: *Pravděpodobnost a matematická statistika*. Skriptum ČVUT, Praha, 2007.
- [Rogalewicz] Rogalewicz, V.: *Pravděpodobnost a statistika pro inženýry*. 2. přepracované vydání, Skriptum FBMI ČVUT, Praha, 2007.
- [Zvára, Štěpán] Zvára, K., Štěpán, J.: *Pravděpodobnost a matematická statistika* (2. vydání). Matfyzpress, MFF UK, Praha, 2002.
- [Anděl: Statistické metody] Anděl, J.: *Statistické metody*. 2. vyd., Matfyzpress, Praha, 1998.
- [Anděl: Matematická statistika] Anděl, J.: *Matematická statistika*. SNTL/Alfa, Praha, 1978.
- [Disman] Disman, M.: *Jak se vyrábí sociologická znalost*. Karolinum, UK, Praha, 2005.
- [Jaroš a kol.] Jaroš, F. a kol.: *Pravděpodobnost a statistika*. Skriptum VŠCHT, 2. vydání, Praha, 1998.
- [Likeš, Machek] Likeš, J., Machek, J.: *Matematická statistika*. 2. vydání, SNTL, Praha, 1988.
- [Nagy] Nagy, I.: *Pravděpodobnost a matematická statistika*. Cvičení. Skriptum FD ČVUT, Praha, 2002.
- [Něničková] Něničková, A.: *Matematická statistika — cvičení*. Skriptum ČVUT, Praha, 1990.
- [Riečanová a kol.] Riečanová, Z. a kol.: *Numerické metody a matematická statistika*. Alfa/SNTL, Bratislava, 1987.
- [Riečan a kol.] Riečan, B., Lamoš, F., Lenárt, C.: *Pravděpodobnost a matematická statistika*. Alfa/SNTL, Bratislava, 1984.
- [SH10] Schlesinger, M.I., Hlaváč, V.: *Deset přednášek z teorie statistického a strukturního rozpoznávání*. ČVUT, Praha, 1999.
- [Swoboda] Swoboda, H.: *Moderní statistika*. Svoboda, Praha, 1977.
- [Chatfield] Chatfield, C.: *Statistics for Technology*. 3rd ed., Chapman & Hall, London, 1992.

- [Hsu] Hsu, H.P.: *Probability, Random Variables, and Random Processes*. McGraw-Hill, 1996.
- [Mood a kol.] Mood, A.M., Graybill, F.A., Boes, D.C.: *Introduction to the Theory of Statistics*. 3rd ed., McGraw-Hill, 1974.
- [Papoulis] Papoulis, A.: *Probability and Statistics*. Prentice-Hall, 1990.
- [Papoulis, Pillai] Papoulis, A., Pillai, S.U.: *Probability, Random Variables, and Stochastic Processes*. 4th ed., McGraw-Hill, Boston, USA, 2002.
- [Spiegel et al. 2000] Spiegel, M.R., Schiller, J.J., Srinivasan, R.A.: *Probability and Statistics*. McGraw-Hill, 2000.
- [Wasserman] Wasserman, L.: *All of Statistics. A Concise Course in Statistical Inference*. Springer, 2004.