

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/stat>

24. listopadu 2020

Obsah

1	O čem to je?	3
1.1	Teorie pravděpodobnosti	3
1.2	Statistika	3
2	Základní pojmy teorie pravděpodobnosti	3
2.1	Náhodný pokus	3
2.2	Laplaceova (klasická) definice pravděpodobnosti	4
2.2.1	Základní pojmy	4
2.3	Vlastnosti pravděpodobnosti	5
2.3.1	Úplný systém jevů	5
2.3.2	Náhodná veličina	5
2.3.3	Střední hodnota	5
2.4	Problémy Laplaceovy definice pravděpodobnosti	6
2.4.1	Rozšíření Laplaceova modelu pravděpodobnosti	6
2.5	Kombinatorické pojmy a vzorce	6
2.5.1	Permutace (pořadí) bez opakování	6
2.5.2	Variace s opakováním (uspořádaný výběr s vrácením)	6
2.5.3	Variace bez opakování (uspořádaný výběr bez vrácení)	6
2.5.4	Kombinace bez opakování (neuspořádaný výběr bez vrácení)	7
2.5.5	Kombinace s opakováním (neuspořádaný výběr s vrácením)	7
2.5.6	Permutace s opakováním	7
2.6	Kolmogorovova definice pravděpodobnosti	8
2.6.1	Borelova σ-algebra	9
2.6.2	Pravděpodobnost (=pravděpodobnostní míra)	9
3	Nezávislost a podmíněná pravděpodobnost	10
3.1	Nezávislé jevy	10
3.2	Podmíněná pravděpodobnost	11
3.2.1	Podmíněná nezávislost	14
4	Náhodné veličiny	16
4.1	Náhodná veličina	16
4.2	Nezávislost náhodných veličin	17
4.3	Směs náhodných veličin	18
4.4	Druhy náhodných veličin	19
4.4.1	Diskrétní náhodné veličiny	19
4.4.2	Spojité náhodné veličiny	19
4.4.3	Smíšené náhodné veličiny	19

4.4.4	Směsi náhodných veličin stejného typu	21
4.5	Kvantilová funkce náhodné veličiny	21
4.6	Jak reprezentovat náhodnou veličinu v počítači	22
4.7	Operace s náhodnými veličinami	23
4.8	Jak realizovat náhodnou veličinu na počítači	26
5	Charakteristiky náhodných veličin	27
5.1	Střední hodnota	27
5.2	Rozptyl (disperze)	28
5.3	Směrodatná odchylka	29
5.4	Obecné a centrální momenty	29
5.5	Normovaná náhodná veličina	29
5.6	Základní typy diskrétních rozdělení	30
5.6.1	Diracova	30
5.6.2	Rovnoměrná	30
5.6.3	Alternativní (Bernoulliho) $\text{Alt}(q)$	31
5.6.4	Binomická $\text{Bi}(m, q)$	31
5.6.5	Poissonova $\text{Po}(\lambda)$	31
5.6.6	Geometrická	32
5.6.7	Hypergeometrická	33
5.7	Základní typy spojitých rozdělení	34
5.7.1	Rovnoměrná $R(a, b)$	34
5.7.2	Normální (Gaussova) $N(\mu, \sigma^2)$	34
5.7.3	Logaritmicko-normální $\text{LN}(\mu, \sigma^2) = \exp(N(\mu, \sigma^2))$	35
5.7.4	Exponenciální $\text{Ex}(\tau)$	35
5.8	Čebyševova nerovnost	36
6	Náhodné vektory	39
6.1	Diskrétní náhodný vektor	40
6.2	Spojité náhodný vektor	40
6.3	Obecnější náhodné veličiny	40
6.4	Číselné charakteristiky náhodného vektoru	41
6.4.1	Vícerozměrné normální rozdělení $N(\mu, \Sigma)$	42
6.5	Reprezentace náhodných vektorů v počítači	42
7	Lineární prostor náhodných veličin	42
7.1	Lineární podprostor \mathcal{N} náhodných veličin s nulovými středními hodnotami	43
7.2	Lineární regrese	43
8	Základní pojmy statistiky	45
8.1	K čemu potřebujeme statistiku	45
8.2	Náhodný výběr, empirické rozdělení	45
8.3	Obecné vlastnosti odhadů	46
8.4	Odhad střední hodnoty	47
8.5	Odhad k -tého obecného momentu EX^k	48
8.6	Odhad rozptylu	49
8.6.1	Odhad rozptylu při známé střední hodnotě	49
8.6.2	Rozdělení χ^2 s n stupni volnosti, $\chi^2(n)$	49
8.6.3	Odhad rozptylu při neznámé střední hodnotě	51
8.6.4	Eficience odhadů rozptylu pro normální rozdělení	53
8.7	Odhad směrodatné odchylky	54
8.8	Histogram a popis empirického rozdělení	54
8.9	Odhad mediánu	55
8.10	Intervalové odhady	55
8.11	Intervalové odhady parametrů normálního rozdělení	55
8.11.1	Intervalový odhad střední hodnoty při známém rozptylu σ^2	55
8.11.2	Intervalový odhad střední hodnoty při neznámém rozptylu	56
8.11.3	Studentovo t-rozdělení	56

8.11.4	Intervalový odhad střední hodnoty při neznámém rozptylu 2	57
8.11.5	Intervalový odhad rozptylu	57
8.11.6	Intervalové odhady spojitých rozdělení, která nejsou normální	58
8.12	Obecné odhady parametrů	59
8.12.1	Metoda momentů	59
8.12.2	Metoda maximální věrohodnosti	60
8.12.3	Příklady na odhady parametrů	62
9	Testování hypotéz	67
9.1	Základní pojmy a principy testování hypotéz	67
9.2	Testy střední hodnoty normálního rozdělení	69
9.2.1	Při známém rozptylu σ^2	69
9.2.2	Při neznámém rozptylu	70
9.3	Testy rozptylu normálního rozdělení	70
9.4	Porovnání dvou normálních rozdělení	70
9.4.1	Testy rozptylu dvou normálních rozdělení [Fisher]	70
9.4.2	Testy středních hodnot dvou normálních rozdělení se stejným známým rozptylem σ^2	72
9.4.3	Testy středních hodnot dvou normálních rozdělení s různými známými rozptily σ_X^2, σ_Y^2	72
9.4.4	Testy středních hodnot dvou normálních rozdělení se stejným neznámým rozptylem σ^2	72
9.4.5	Testy středních hodnot dvou normálních rozdělení - párový test	73
9.5	Korelace, její odhad a testování	74
9.5.1	Test nekorelovanosti dvou normálních rozdělení	74
9.6	χ^2 -test dobré shody	74
9.6.1	Základní podoba testu	74
9.6.2	Modifikace	76
9.6.3	χ^2 -test nezávislosti dvou rozdělení	76
9.6.4	χ^2 -test dobré shody dvou rozdělení	77
9.7	Neparametrické testy	77
9.7.1	Znaménkový test	77
9.7.2	Wilcoxonův test (jednovýběrový)	78

1 O čem to je?

1.1 Teorie pravděpodobnosti

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

Poskytuje model a kvantifikaci výsledků.

Pravděpodobnostní **popis** \Rightarrow **chování** systému

1.2 Statistika

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

Chování systému \Rightarrow pravděpodobnostní **popis**

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

2 Základní pojmy teorie pravděpodobnosti

2.1 Náhodný pokus

„Takový, na který si můžeme vsadit.“

Tedy nikoli:

- Jak je pravděpodobné, že ve skriptech na str. 42 je chyba?
- Jak je pravděpodobné, že zítra bude v menze dobrý oběd?

Realizace náhodného pokusu – např. losovací zařízení:

- Kostka, čtyřstěn, dvanáctistěn...
- Tužka, dlouhý hranol...
- „Kolo štěstí.“
- Urna s losy.

2.2 Laplaceova (klasická) definice pravděpodobnosti

Předpoklad: Náhodný pokus s $n \in \mathbb{N}$ různými, po dvou neslučitelnými výsledky, které jsou **stejně možné**. Jev, který nastává právě při k z těchto výsledků, má pravděpodobnost k/n .

(Urna s n losy, z nichž k „vyhrává“.)

1. problém: Co to je „stejně možné“? „Stejně pravděpodobné“? (definice kruhem!)

Elementární jevy jsou všechny „stejně možné“ výsledky (*losy*).

Množina všech elementárních jevů: Ω (*urna*)

Jev: $A \subseteq \Omega$ (*množina vyhrávajících losů*)

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

2.2.1 Základní pojmy

Jev jistý: Ω , **1** (*všechny losy vyhrávají*)

Jev nemožný: \emptyset , **0** (*žádný los nevyhrává*)

Konjunkce jevů („and“): $A \cap B$ (*losy, které vyhrávají v obou tazích*)

Disjunkce jevů („or“): $A \cup B$ (*losy, které vyhrávají v aspoň jednom tahu*)

Jev opačný k A : $\bar{A} = \Omega \setminus A$ (*losy, které nevyhrávají*)

$A \Rightarrow B$: $A \subseteq B$

Jevy neslučitelné: $A_1, \dots, A_n : \bigcap_{i \leq n} A_i = \emptyset$

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 $\exp \Omega$ (prozatím množina všech podmnožin množiny Ω)

Pravděpodobnost jevu A :

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

kde $|\cdot|$ značí počet prvků množiny

2.3 Vlastnosti pravděpodobnosti

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

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

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

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

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

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

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

2.3.1 Úplný systém jevů

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

Příklad: Ruleta: např.

A. $\{0\}, \{1\}, \dots, \{36\}$;

B. sudá, lichá;

C. $\{0\}$, červená, černá.

Speciální případ pro 2 jevy: $\{C, \bar{C}\}$, $B_1 = C, B_2 = \bar{C}$.

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

Motivační příklad (kolik je nekuřáků): Mužů je v populaci 48 %, kuřáků a kuřáček dohromady 30 %. Jakých hodnot může nabývat pravděpodobnost, že náhodně vybraný člověk je muž a nekuřák?

M ... muž, $P(M) = 0.48$

K ... kuřák, $P(K) = 0.3, P(\bar{K}) = 0.7$

Hledaná pravděpodobnost $P(M \cap \bar{K})$ jevu opačného k $\bar{M} \cup K$,

$$\max\{P(\bar{M}), P(K)\} = 0.52 \leq P(\bar{M} \cup K) \leq P(\bar{M}) + P(K) = 0.52 + 0.3 = 0.82$$

$$1 - 0.82 = 0.18 \leq P(M \cap \bar{K}) \leq 1 - 0.52 = 0.48$$

Možné jsou všechny hodnoty z intervalu $\langle 0.18, 0.48 \rangle$.

2.3.2 Náhodná veličina

je (prozatím libovolná) **funkce** $X: \Omega \rightarrow \mathbb{R}$

Příklady: Teplota v Klementinu zítra v 6 hodin.

Známka ze zkoušky.

Kolo štěstí. Na každém poli je napsána částka výhry.

Toto je univerzální příklad na pravděpodobnost dle Laplace.

2.3.3 Střední hodnota

Příklad: Kolo štěstí má 4 stejně velké sekce (A, B, C, D) , na které připadají výhry 0, 10, 50, 100.

Elementární jevy jsou všechny možné výsledky,

$$\Omega = \{A, B, C, D\},$$

náhodná veličina X je výše výhry:

ω	A	B	C	D
$X(\omega)$	0	10	50	100

Střední hodnota je spravedlivá cena za účast ve hře:

$$EX = \frac{0 + 10 + 50 + 100}{4} = \frac{160}{4} = 40.$$

Obecně

$$EX = \frac{1}{|\Omega|} \sum_{\omega \in \Omega} X(\omega).$$

2.4 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 prováděn „rovnoměrně“).

3. problém: Nedovoluje iracionální hodnoty pravděpodobnosti.

Příklad: Kolo štěstí s nestejnými sekcemi.

2.4.1 Rozšíření Laplaceova modelu pravděpodobnosti

Elementární jevy nemusí být stejně pravděpodobné.

Ztrácíme návod, jak vybrat „správnou“ pravděpodobnost.

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.

To je role statistiky, která k danému opakovatelnému pokusu hledá pravděpodobnostní model.

2.5 Kombinatorické pojmy a vzorce

(Dle [Zvára, Štěpán].)

Losovací zařízení:

Urna s losy, které nelze rozlišit před losováním a lze rozlišit po vylosování.

Losů je n (zde 9, číslovaných 1, 2, 3, 4, 5, 6, 7, 8, 9).

Větší počet výsledků umožní **opakované losování**, které lze provést různě.

2.5.1 Permutace (pořadí) bez opakování

Vytáhneme všechny losy, záleží na pořadí.

1.	2.	3.	4.	5.	6.	7.	8.	9.
2	7	1	8	3	9	4	5	6

$$9 \cdot 8 \cdot 7 \cdot 6 \cdot 5 \cdot 4 \cdot 3 \cdot 2 \cdot 1 = 9!$$

Obecně $n!$ výsledků stejně pravděpodobných

Nadále postupně vytáhneme k losů (zde 6).

2.5.2 Variace s opakováním (uspořádaný výběr s vrácením)

(záleží na pořadí, např. číselný zámek)

1.	2.	3.	4.	5.	6.
2	7	1	8	2	8

9^6

Obecně n^k výsledků stejně pravděpodobných

2.5.3 Variace bez opakování (uspořádaný výběr bez vrácení)

(záleží na pořadí, $k \leq n$)

1.	2.	3.	4.	5.	6.
2	7	1	8	23	89

$$9 \cdot 8 \cdot 7 \cdot 6 \cdot 5 \cdot 4 = \frac{9!}{3!}$$

Obecně $\frac{n!}{(n-k)!}$ výsledků stejně pravděpodobných

2.5.4 Kombinace bez opakování (neuspořádaný výběr bez vracení)

(nezáleží na pořadí, $k \leq n$)

1	2	3	4	5	6	7	8	9
x	x	x				x	x	x

$$\frac{9 \cdot 8 \cdot 7 \cdot 6 \cdot 5 \cdot 4}{6 \cdot 5 \cdot 4 \cdot 3 \cdot 2 \cdot 1} = \frac{9!}{3!6!} = \binom{9}{6}$$

Obecně $\binom{n}{k}$ výsledků stejně pravděpodobných

2.5.5 Kombinace s opakováním (neuspořádaný výběr s vracením)

(nezáleží na pořadí)

1	2	3	4	5	6	7	8	9
x	x					x	x	
	x						x	

$$\begin{array}{c} 1 \mid 22 \mid \mid \mid \mid 7 \mid 88 \mid \\ \cdot \mid \cdot \mid \mid \mid \mid \cdot \mid \cdot \mid \end{array}$$

$$\binom{14}{6}$$

Obecně $\binom{n+k-1}{k}$ výsledků **nestějně** pravděpodobných.

1	2	3	4	5	6	7	8	9
x								
x								
x								
x								
x								
x								

$$P(111111 \mid \mid \mid \mid \mid \mid) = \frac{1}{9^6}$$

\ll

1	2	3	4	5	6	7	8	9
x	x	x	x	x	x			

$$P(1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid \mid \mid) = \frac{6!}{9^6}$$

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

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

2.5.6 Permutace s opakováním

Některé losy mohou být **nerozlišitelné**, např. 1, 1, 1, 2, 2, 3, 3, 4, 5,

obecněji $\underbrace{1, \dots, 1}_{k_1 \times}, \underbrace{2, \dots, 2}_{k_2 \times}, \dots, \underbrace{i, \dots, i}_{k_i \times}$, kde $\sum_{j=1}^i k_j = n$.

(záleží na pořadí)

1.	2.	3.	4.	5.	6.	7.	8.	9.
3	1	4	1	5	3	2	1	2

$3! \cdot 2! \cdot 2!$ (obecně $k_1! \cdot \dots \cdot k_i!$) permutací dává stejný výsledek.

$\frac{9!}{3! \cdot 2! \cdot 2!}$ (obecně $\frac{n!}{k_1! \cdot \dots \cdot k_i!}$) výsledků stejně pravděpodobných.

Speciálně pro $i = 2$: $\frac{n!}{k_1! \cdot k_2!} = \frac{n!}{k_1! \cdot (n - k_1)!} = \binom{n}{k_1}$

(k_1 -prvkové kombinace bez opakování z n prvků;

rozlišujeme pouze k_1 „vylosovaných“ a $n - k_1$ „nevylosovaných“).

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. 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.

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

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. Pro $n \gg k$ je počet variací (resp. kombinací) bez opakování přibližně

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

Jednodušší bývají **variace s opakováním** (uspořádaný výběr s vrácením) nebo **kombinace bez opakování** (neuspořádaný výběr bez vrácení).

2.6 Kolmogorova definice pravděpodobnosti

Elementární jevy = všechny možné výsledky pokusu = prvky množiny Ω .

Může jich 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) $\emptyset \in \mathcal{A}$.

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

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

Systém \mathcal{A} podmnožin nějaké množiny Ω , který splňuje podmínky (A1)–(A3), se nazývá **σ -algebra**.
Důsledky: $\Omega = \emptyset \in \mathcal{A}$,

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

Maximalista: Volme $\mathcal{A} = \exp \Omega$.

Vede k nežádoucím problémům, např. Banachův-Tarského paradox.

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

Maximalista: Volme uzavřenost na **jakákoli** sjednocení.

A jsme tam, kde jsme byli...

Praktik: Volme uzavřenost na **konečná** sjednocení.

Nedovoluje např. vyjádřit kruh jako sjednocení obdélníků.

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

2.6.1 Borelova σ -algebra

$\mathcal{B}(\mathbb{R})$ je nejmenší σ -algebra podmnožin \mathbb{R} , která obsahuje všechny intervaly.

Obsahuje všechny intervaly otevřené, uzavřené, polouzavřené, jejich spočetná sjednocení, některé další množiny (např. Cantorovo diskontinuum), ale ne všechny. Její prvky nazýváme **borelovské množiny**.

2.6.2 Pravděpodobnost (=pravděpodobnostní míra)

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

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

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

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, se nazývá **pravděpodobnostní prostor**.

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

Praktik: Spokojme se s **konečnou** aditivitou.

Problémem je např. pravděpodobnost výsledku v kruhu coby sjednocení obdélníků (spočetně mnoha).

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

Maximalista: Požadujeme **úplnou** aditivitu (pro jakékoli soubory po dvou neslučitelných jevů).

Pak bychom nepřipouštěli žádné spojitě rozdělení (ani rovnoměrné ani normální).

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 i iracionální
$P(A) = 0 \Rightarrow A = \emptyset$	možné jevy s nulovou p-stí
p-sti určeny strukturou jevů	p-sti neurčeny strukturou jevů

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

$$\frac{2}{\pi} \doteq 0.63661977236758134307553505349005744.$$

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 $P(A \cap B) = 0$ 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:

Příklad. Máme dva hody mincí a jevy

A ... při prvním hodu padne líc,

B ... při druhém hodu padne líc,

C ... při právě jednom hodu padne líc.

$$\begin{aligned} P(A) &= P(B) = P(C) = \frac{1}{2}, \\ P(A \cap B) &= P(A \cap C) = P(B \cap C) = \frac{1}{4} = P(A) \cdot P(B) = P(A) \cdot P(C) = P(B) \cdot P(C), \\ P(A \cap B \cap C) &= 0 \neq \frac{1}{8} = P(A) \cdot P(B) \cdot P(C), \end{aligned}$$

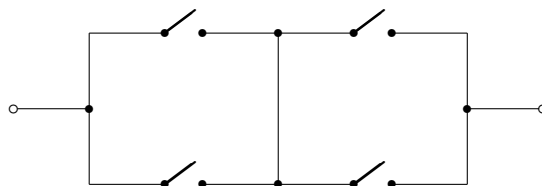
takže jevy A, B, C jsou po dvou nezávislé, ale nejsou nezávislé.

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 $\mathcal{K} \subseteq \mathcal{M}$.

Příklad. Spínače zapojené dle obrázku jsou nezávisle sepnuty s pravděpodobností 0.9. S jakou pravděpodobností celá soustava povede proud?



Dvě sériově spojené části, každá vede s pravděpodobností

$$0.9 + 0.9 - 0.9^2 = 0.99,$$

celek vede s pravděpodobností

$$0.99^2 = 0.9801.$$

Příklad. (Sally Clark): Syndrom náhlého úmrtí kojenců se vyskytuje s pravděpodobností $1/8500$. Jaká je pravděpodobnost dvou úmrtí po sobě v jedné rodině?

Soud uvěřil žalobci, že $1/8500^2 = 1/72\,250\,000$, a matku odsoudil.

Později osvobozena ona i 3 další ženy.

I kdyby platila nezávislost, velmi pravděpodobně by byl někdo neprávem odsouzen.

3.2 Podmíněná pravděpodobnost

Motivační příklad (alkohol za volantem):

90 % všech nehod způsobili střízliví řidiči.

Alkoholik: Když se napiju, budu mít $9\times$ menší riziko havárie.

Statistik: To by byla pravda, kdyby opilých bylo stejně jako střízlivých.

Příklad: Pravděpodobnosti výsledků tenisového zápasu se podstatně změnil po odehrání prvního setu.

Máme pravděpodobnostní popis systému. Dostaneme-li dodatečnou informaci, že nastal jev B , můžeme aktualizovat naši znalost o pravděpodobnosti libovolného jevu A . Ten lze vyjádřit jako **disjunktní** sjednocení $(A \cap B) \cup (A \cap \bar{B})$, takže

$$P(A) = P(A \cap B) + P(A \cap \bar{B}).$$

Je-li $P(B) \neq 0 \neq P(\bar{B})$, můžeme roznásobit:

$$P(A) = P(B) \underbrace{\frac{P(A \cap B)}{P(B)}}_{P(A|B)} + P(\bar{B}) \underbrace{\frac{P(A \cap \bar{B})}{P(\bar{B})}}_{P(A|\bar{B})}.$$

Funkce $P(\cdot|B), P(\cdot|\bar{B}): \mathcal{A} \rightarrow \langle 0, 1 \rangle$,

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

jsou pravděpodobnosti na \mathcal{A} , neboť splňují

$$(P1) \quad P(\mathbf{1}|B) = \frac{P(\mathbf{1} \cap B)}{P(B)} = \frac{P(B)}{P(B)} = 1$$

a pro $A_n, n \in \mathbb{N}$, po dvou neslučitelné

$$(P2) \quad P\left(\bigcup_{n \in \mathbb{N}} A_n \middle| B\right) = \frac{P\left(\left(\bigcup_{n \in \mathbb{N}} A_n\right) \cap B\right)}{P(B)} = \frac{P\left(\bigcup_{n \in \mathbb{N}} (A_n \cap B)\right)}{P(B)} = \frac{\sum_{n \in \mathbb{N}} P(A_n \cap B)}{P(B)} = \\ = \sum_{n \in \mathbb{N}} P(A_n|B).$$

(Obdobně pro $P(\cdot|\bar{B})$.)

Nazývají se **podmíněné pravděpodobnosti**.

$P(A|B)$ čteme např. „pravděpodobnost jevu A za podmínky B .“

Je-li $P(A|B)$ definována, jsou jevy A, B **nezávislé**, právě když $P(A|B) = P(A)$.

Podmíněné pravděpodobnosti navíc splňují $B \subseteq A \Rightarrow P(A|B) = 1$, $P(A \cap B) = 0 \Rightarrow P(A|B) = 0$, speciálně $P(B|B) = 1$, $P(\bar{B}|B) = 0$.

(Obdobně pro $P(\cdot|\bar{B})$.)

Původní pravděpodobnost $P(\cdot)$ jsme vyjádřili jako konvexní kombinaci pravděpodobností $P(\cdot|B)$, $P(\cdot|\bar{B})$, odpovídajících situacím, kdy jev B nastal, resp. nenastal:

$$P(A) = P(B) P(A|B) + P(\bar{B}) P(A|\bar{B}).$$

Tato podmínka spolu s $P(B|B) = 1 = P(\bar{B}|\bar{B})$ určuje pravděpodobnosti $P(\cdot|B)$, $P(\cdot|\bar{B})$ jednoznačně. (Pokud není jedna z pravděpodobností $P(B)$, $P(\bar{B})$ nulová.)

Obecněji:

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:

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

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):

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

Motivační příklad (test na drogy): Policií užívaný test DrugWipe 5S je falešně pozitivní u 5 % nezdrogovaných a falešně negativní u 3 % zdrogovaných. Jaká je (podmíněná) pravděpodobnost, že řidič s pozitivním testem je pod vlivem drog?

D ... zdrogovaný, T ... pozitivní test,

$$P(T|\bar{D}) = 0.05, \quad P(T|D) = 1 - P(\bar{T}|D) = 1 - 0.03 = 0.97.$$

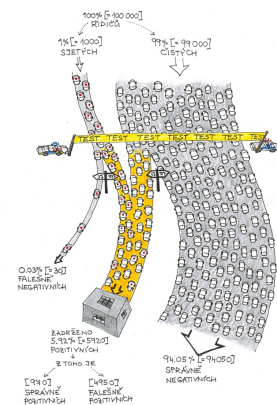
Potřebujeme ještě znát podíl zdrogovaných řidičů na silnicích $P(D)$.

Policejní ani statistická data nemáme, tedy jen zkusíme, co bychom dostali pro

$$P(D) = 0.01, \quad P(\bar{D}) = 0.99:$$

$$\begin{aligned} P(T) &= P(T|D) \cdot P(D) + P(T|\bar{D}) \cdot P(\bar{D}) \\ &= 0.97 \cdot 0.01 + 0.05 \cdot 0.99 = 0.0097 + 0.0495 = 0.0592 \end{aligned}$$

$$\begin{aligned} P(D|T) &= \frac{P(D) \cdot P(T|D)}{P(D) \cdot P(T|D) + P(\bar{D}) \cdot P(T|\bar{D})} \\ &= \frac{P(D) \cdot P(T|D)}{P(T)} = \frac{0.0097}{0.0592} \doteq 0.164. \end{aligned}$$



Zdroj: <https://finmag.penize.cz/kaleidoskop/407932-velka-drogova-kocovina>

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 jestliže jsme pozorovali A .

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

Příklad: Informační kanál

B_j ... vyslán j -tý vstupní znak, $j \in \{1, \dots, m\}$

A_i ... přijat i -tý výstupní znak, $i \in \{1, \dots, k\}$ (může být $k \neq m$)

Lze odhadnout podmíněné pravděpodobnosti $P(A_i|B_j)$, že znak j bude přijat jako i .

Z apriorních pravděpodobností (vyslání znaku j) $P(B_j)$ můžeme maticovým násobením určit pravděpodobnosti přijatých znaků:

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

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

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

Rozdělení pravděpodobností vyslaných znaků je

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

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

Motivační příklad (alkohol za volantem – pokračování):

90 % všech nehod způsobili střízliví řidiči.

99 % řidičů bylo střízlivých.

Označme jevy

A ... požil alkohol, $P(A) = 0.01$,

H ... způsobil nehodu, $P(A|H) = 0.1$.

$$\begin{aligned}
0.1 &= P(A|H) = \frac{P(A) \cdot P(H|A)}{P(A) \cdot P(H|A) + P(\bar{A}) \cdot P(H|\bar{A})} = \\
&= \frac{0.01 \cdot P(H|A)}{0.01 \cdot P(H|A) + 0.99 \cdot P(H|\bar{A})} = \frac{1}{1 + 99 \cdot \frac{P(H|\bar{A})}{P(H|A)}}.
\end{aligned}$$

Požitím alkoholu se zvyšuje riziko nehody

$$\frac{P(H|A)}{P(H|\bar{A})} = 11 \times .$$

Kdyby bylo 50 % řidičů opilých, $P(A) = 0.5$, jejich podíl na haváriích by byl

$$\begin{aligned}
P(A|H) &= \frac{P(A) \cdot P(H|A)}{P(A) \cdot P(H|A) + P(\bar{A}) \cdot P(H|\bar{A})} = \\
&= \frac{0.5 \cdot P(H|A)}{0.5 \cdot P(H|A) + 0.5 \cdot P(H|\bar{A})} = \frac{1}{1 + \frac{P(H|\bar{A})}{P(H|A)}} = \frac{1}{1 + \frac{1}{11}} = \frac{11}{12}.
\end{aligned}$$

(Neuvažovali jsme, že účastníků nehody bývá víc a přítomnost alkoholu u jejích účastníků nemusí být nezávislá.)

Příklad: Detektor v CERNu zachytí *událost* (současný vznik Higgsova bosonu a páru top kvarků) s účinností 0.01;

falešně pozitivní je s pravděpodobností 0.000 05.

Hledaná událost se vyskytne v 6 případech z 10^{12} . Jaká je pravděpodobnost, že detekovaná událost skutečně nastala? ¹

Označme jevy

U ... událost nastala, $P(U) = 6 \cdot 10^{-12}$,

D ... detektor hlásí událost, $P(D|U) = 0.01$, $P(D|\bar{U}) = 5 \cdot 10^{-5}$.

$$\begin{aligned}
P(D) &= P(U) \cdot P(D|U) + P(\bar{U}) \cdot P(D|\bar{U}) = \\
&= 6 \cdot 10^{-12} \cdot 0.01 + (1 - 6 \cdot 10^{-12}) \cdot 5 \cdot 10^{-5} = 5.000000006 \cdot 10^{-5},
\end{aligned}$$

$$P(U|D) = \frac{P(U) \cdot P(D|U)}{P(D)} = \frac{6 \cdot 10^{-14}}{5 \cdot 10^{-5}} = 1.2 \cdot 10^{-9}.$$

Příklad. (Sally Clark — pokračování): Jaká je (podmíněná) pravděpodobnost syndromu náhlého úmrtí kojenců v rodině, kde se už vyskytl?

Není dostatek dat pro kvalifikovaný odhad.

Alespoň jsou známy vlivy, které (podmíněné) riziko snižují (nekuřácká rodina, stabilní, dobře zajištěná).

Tytéž vlivy snižují i riziko vraždy novorozence, takže **podmíněná** pravděpodobnost, že šlo o trestný čin, zůstává nízká.

3.2.1 Podmíněná nezávislost

Příklad. A. Kolik reálných parametrů (stupňů volnosti) je třeba pro určení pravděpodobností všech jevů, které lze popsat logickými výrazy z n výchozích jevů?

B. Jak se toto číslo změní, předpokládáme-li, že výchozí jevy jsou nezávislé?

Řešení. A. Každý z uvažovaných jevů lze vyjádřit v úplné disjunktivní normální formě jako disjunkci výrazů, které jsou konjunkcemi n výchozích jevů nebo jejich negací; těch je 2^n a tvoří úplný systém jevů, popsaných $2^n - 1$ parametry (součet jejich pravděpodobností 1).

B. Stačí n pravděpodobností výchozích jevů.

¹André Sopczak: Recognition and categorization of events resulting from proton collisions at CERN. Přednáška na katedře kybernetiky FEL ČVUT, 22. 1. 2019

n	2	3	4	5	10	16	30
$2^n - 1$	3	7	15	31	1023	65 535	1 073 741 823

Je vidět, že předpoklad nezávislosti může vést na mnohem jednodušší model. Často je však neopodstatněný.

Motivační příklad (popis systému se závislými jevy): Chceme popsat pravděpodobnosti jevů, složených z následujících 4 charakteristik lidí:

- rád nosí růžové oblečení,
- používá make-up,
- pracuje jako zdravotní sestra,
- prodělal rakovinu prsu.

Byla zjištěna velmi podstatná závislost; potřebujeme $2^4 - 1 = 15$ parametrů.

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

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

Obdobně definujeme podmíněnou nezávislost více jevů.

Motivační příklad (popis systému se závislými jevy) – řešení: Přidáme ještě pátý parametr:

- je žena,

a shledáme, že při jeho znalosti lze ostatní 4 parametry považovat za podmíněně nezávislé. Stačí nám 4 parametry pro ženy, 4 pro muže a 1 pro pravděpodobnost, že osoba je žena, tj. 9 parametrů místo 15, resp. 31.

4 Náhodné veličiny

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

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 nesetkáme.

Pravděpodobnostní míra P_X splňuje podmínky:

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

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

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

$$I \subseteq J \Rightarrow P_X(I) \leq P_X(J), \quad P_X(J \setminus I) = P_X(J) - P_X(I).$$

Popisy náhodné veličiny

prostor elementárních jevů	Ω	\mathbb{R}
σ -algebra jevů	\mathcal{A}	$\mathcal{B}(\mathbb{R})$
pravděpodobnostní míra	P	P_X
pravděpodobnostní prostor	(Ω, \mathcal{A}, P)	$(\mathbb{R}, \mathcal{B}(\mathbb{R}), P_X)$
náhodná proměnná	$X: \Omega \rightarrow \mathbb{R}, \quad \omega \mapsto X(\omega)$	$\text{id}: \mathbb{R} \rightarrow \mathbb{R}, \quad x \mapsto x$
$P(X \in I)$	$P(\{\omega \in \Omega \mid X(\omega) \in I\})$	$P_X(I)$

Příklad: Počet figurek Člověče nezlob se!, které vstupují do hry po jednom hodu kostkou.

prostor elementárních jevů	$\Omega = \{1, 2, 3, 4, 5, 6\}$	\mathbb{R}
σ -algebra jevů	$\mathcal{A} = \exp \Omega$	$\mathcal{B}(\mathbb{R})$
pravděpodobnostní míra	$P(A) = \frac{ A }{6}$	$P_X(I) = \begin{cases} 1, & 0, 1 \in I \\ 5/6, & 0 \in I, 1 \notin I \\ 1/6, & 0 \notin I, 1 \in I \\ 0, & 0, 1 \notin I \end{cases}$
náhodná proměnná	$X(\omega) = \begin{cases} 1, & \omega = 6 \\ 0, & \text{jinak} \end{cases}$	$X: x \mapsto x$

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

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

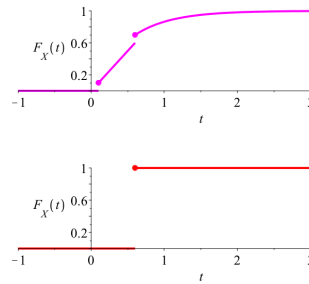
$F_X: \mathbb{R} \rightarrow \langle 0, 1 \rangle$ je **distribuční funkce** náhodné veličiny X . Ta stačí, neboť

$$\begin{aligned} (a, b) &= (-\infty, b) \setminus (-\infty, a), & P_X((a, b)) &= P(a < X \leq b) = F_X(b) - F_X(a), \\ (a, \infty) &= \mathbb{R} \setminus (-\infty, a), & P_X((b, \infty)) &= 1 - F_X(b), \\ \{b\} &= \bigcap_{n \rightarrow \infty} (b - \frac{1}{n}, b) & P_X(\{a\}) &= P(X = a) = \lim_{n \rightarrow \infty} (F_X(b) - F_X(b - \frac{1}{n})) \\ &= \lim_{n \rightarrow \infty} (b - \frac{1}{n}, b) & &= F_X(b) - \lim_{a \rightarrow b-} F_X(a), \\ \dots & & \dots & \end{aligned}$$

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: $X \leq Y \Rightarrow F_X \geq F_Y$,
neboli $(\forall \omega \in \Omega : X(\omega) \leq Y(\omega)) \Rightarrow (\forall t \in \mathbb{R} : F_X(t) \geq F_Y(t))$.

Důkaz.

$$\begin{aligned} \forall t \in \mathbb{R} : (Y(\omega) \leq t \Rightarrow X(\omega) \leq Y(\omega) \leq t), \\ \{\omega \in \Omega : Y(\omega) \leq t\} \subseteq \{\omega \in \Omega : X(\omega) \leq t\}, \\ F_Y(t) = P(Y \leq t) = P(\{\omega \in \Omega : Y(\omega) \leq t\}) \leq P(\{\omega \in \Omega : X(\omega) \leq t\}) = P(X \leq t) = F_X(t). \end{aligned}$$

□

Příklad: V zimním období je venkovní teplota, X , nižší než vnitřní, Y . Např. pro $t = 0$ tvrdíme, že pravděpodobnost, že mrzne uvnitř, je menší než pravděpodobnost, že mrzne venku.

4.2 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.

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

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

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

Ekvivalentně stačí požadovat

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

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

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 $I_i = \mathbb{R}$, resp. $t_i = \infty$. Pak $P(X_i \in I_i) = 1$ a v součinu se tento činitel neprojeví.

Spočetná nekonečná množina náhodných veličin je nezávislá, je-li každá její konečná podmnožina nezávislá.

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 .

4.3 Směs náhodných veličin

Příklad: Náhodné veličiny V, U jsou výsledky studenta při odpovědích na dvě zkouškové 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 V , resp. U 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

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

$P(A_1 \cup A_2) = c P_1(A_1) + (1 - c) P_2(A_2)$ pro $A_1 \in \mathcal{A}_1, A_2 \in \mathcal{A}_2$.

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

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

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

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

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

a distribuční funkci

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

Podobně definujeme obecněji **směs náhodných veličin** V_1, \dots, V_n s **koefficienty** $c_1, \dots, c_n \in \langle 0, 1 \rangle$, $\sum_{i=1}^n c_i = 1$, značíme $\text{Mix}_{(c_1, \dots, c_n)}(V_1, \dots, V_n) = \text{Mix}_{\mathbf{c}}(V_1, \dots, V_n)$, kde $\mathbf{c} = (c_1, \dots, c_n)$. Má pravděpodobnostní míru $\sum_{i=1}^n c_i P_{V_i}$ a distribuční funkci $\sum_{i=1}^n c_i F_{V_i}$. (Lze zobecnit i na spočetně mnoho náhodných veličin.)

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

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

Příklad: Směsí reálných čísel r_1, \dots, r_n s koefficienty 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.

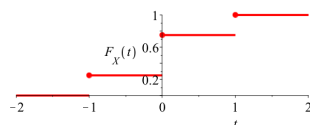
4.4 Druhy náhodných veličin

4.4.1 Diskrétní náhodné veličiny

(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

$$\Omega_X = \{t \in \mathbb{R} : P_X(\{t\}) \neq 0\} = \{t \in \mathbb{R} : P(X = t) \neq 0\}.$$

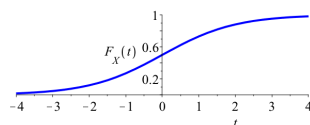


Diskrétní náhodnou veličinu popisuje **pravděpodobnostní funkce** $p_X(t) = P_X(\{t\}) = P(X = t)$. Splňuje

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

4.4.2 Spojité náhodné veličiny

Mají spojitou distribuční funkci.



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 .

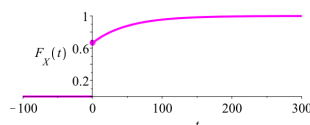
Lze volit $f_X(t) = \frac{dF_X(t)}{dt}$, pokud derivace existuje.

$P_X(\{t\}) = 0$ pro všechna t .

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

4.4.3 Smíšené náhodné veličiny

Motivační příklad: Náhodné veličiny jako množství srážek, spotřeba elektřiny apod. je vhodné modelovat spojitým rozdělením, jsou-li nenulové, nicméně nulové mohou být s nezanedbatelnou nenulovou pravděpodobností.



Podobně i u dalších náhodných veličin se mohou uplatnit meze či jiné preferované hodnoty, nastávající s nenulovou pravděpodobností.

Směs předchozích dvou případů;

$\Omega_X \neq \emptyset$, $P_X(\mathbb{R} \setminus \Omega_X) = P(X \notin \Omega_X) \neq 0$.

Nelze je popsat ani pravděpodobnostní funkcí (existuje, ale neurčuje celé rozdělení) ani hustotou (neexistuje, nevychází konečná).

Každou náhodnou veličinu se smíšeným rozdělením lze **jednoznačně** vyjádřit ve tvaru $X = \text{Mix}_c(V, U)$, kde V je spojitá, U je diskrétní a $c \in (0, 1)$:

Nespojitostí je spočetně mnoho, lze je očíslovat; n -tá je v bodě r_n a má velikost

$$c_n := F_X(r_n) - \lim_{t \rightarrow r_n^-} F_X(t).$$

Odpovídá jí složka směsi r_n s Diracovým rozdělením a váhou c_n .

$$F_X(t) = \sum_n c_n F_{r_n}(t) + G(t),$$

$$G := F_X - \sum_n c_n F_{r_n}$$

je spojitá neklesající funkce,

$$\lim_{t \rightarrow -\infty} G(t) = 0,$$

$$\lim_{t \rightarrow \infty} G(t) = 1 - \sum_n c_n =: c,$$

$$F_V := \frac{G}{c}$$

je distribuční funkce spojitě náhodné veličiny V ,

$$X = \text{Mix}_{(c, c_1, c_2, \dots)}(V, r_1, r_2, \dots),$$

$$U = \underbrace{\frac{1}{\sum_n c_n}}_{1-c} \text{Mix}_{(c_1, c_2, \dots)}(r_1, r_2, \dots),$$

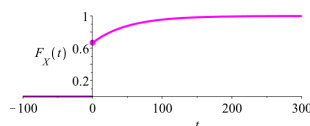
je diskrétní složka náhodné veličiny X ,

$$X = \text{Mix}_{(c, 1-c)}(V, U).$$

Motivační příklad (dešťové srážky):

Srážkový úhrn v mm za 24 hodin má rozdělení s distribuční funkcí

$$F_X(t) = \begin{cases} 1 - \frac{1}{3} \exp(-\frac{t}{50}), & t \geq 0, \\ 0 & \text{jinak.} \end{cases}$$



Po 2/3 dní neprší, diskrétní složka je $U = 0$ s váhou $c_1 := 2/3$.

Motivační příklad (dešťové srážky):

Srážkový úhrn v mm za 24 hodin má rozdělení s distribuční funkcí

$$F_X(t) = \begin{cases} 1 - \frac{1}{3} \exp(-\frac{t}{50}), & t \geq 0, \\ 0 & \text{jinak.} \end{cases}$$

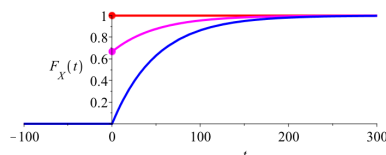
Po $2/3$ dní neprší, diskrétní složka je $U = 0$ s váhou $c_1 := 2/3$.
 Spojitá složka V má váhu $c := 1/3$, distribuční funkci

$$F_V(t) = \begin{cases} 1 - \exp(-\frac{t}{50}), & t \geq 0, \\ 0 & \text{jinak} \end{cases}$$

a hustotu

$$f_V(t) = \begin{cases} \frac{1}{50} \exp(-\frac{t}{50}), & t \geq 0, \\ 0 & \text{jinak} \end{cases}$$

(exponenciální rozdělení).



4.4.4 Směsi náhodných veličin stejného typu

$$X = \text{Mix}_{(c,1-c)}(V, U).$$

Jsou-li V, U **diskrétní**, má X pravděpodobnostní funkci

$$p_X = c p_V + (1 - c) p_U.$$

Jsou-li V, U **absolutně spojité**, má X hustotu

$$f_X = c f_V + (1 - c) f_U.$$

Obdobně pro směsi více náhodných veličin.

4.5 Kvantilová funkce náhodné veličiny

Příklad. 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 $F_X(t) = \alpha$. To nemusí existovat, ale vždy existuje t , pro které

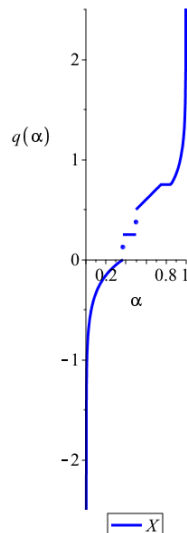
$$P(X < t) \leq \alpha \leq P(X \leq t),$$

tj.

$$\lim_{u \rightarrow t-} F_X(u) \leq \alpha \leq F_X(t),$$

Všechna taková t tvoří omezený interval, z něhož bereme (obvykle) střed,

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



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,
- $q_X(\alpha) = \frac{1}{2} \left(\lim_{\beta \rightarrow \alpha-} q_X(\beta) + \lim_{\beta \rightarrow \alpha+} q_X(\beta) \right)$.

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).

4.6 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 $2n$ čísly (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.

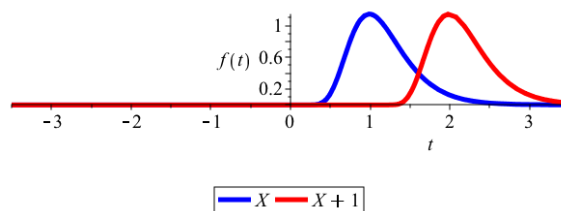
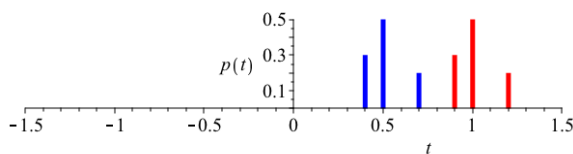
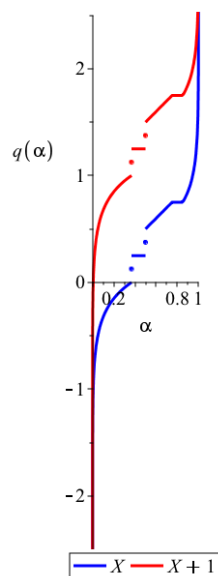
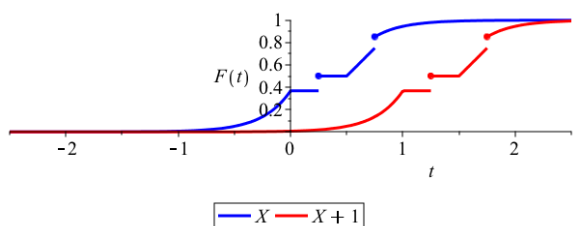
3. **Smíšená:** Jako u spojitě náhodné veličiny. Tento popis je však pro diskrétní část zbytečně nepřesný. Můžeme použít rozklad na diskrétní a spojitou část.

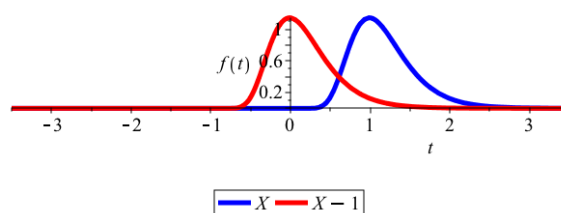
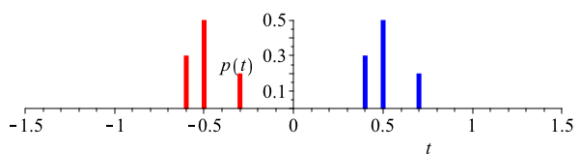
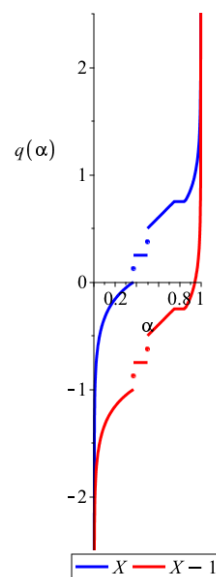
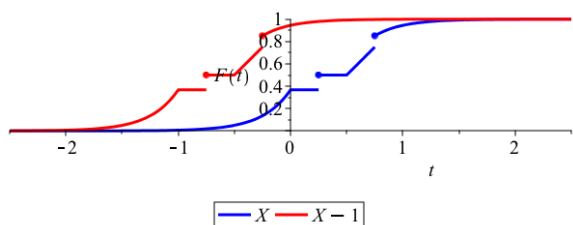
4.7 Operace s náhodnými veličinami

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

Přičtení konstanty r odpovídá posunutí ve směru vodorovné osy:

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



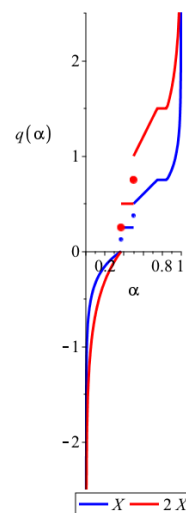
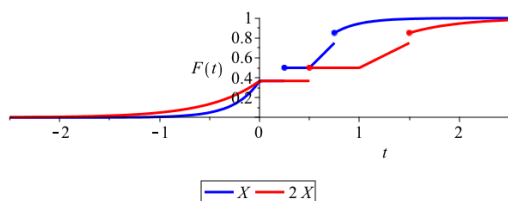


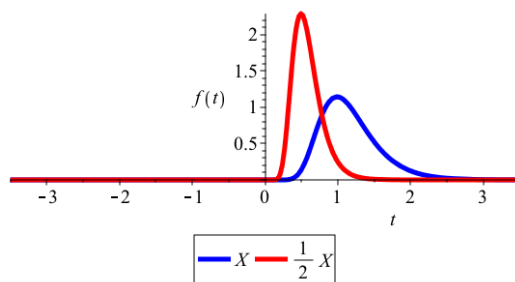
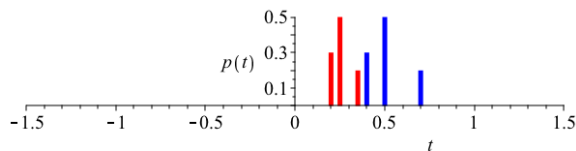
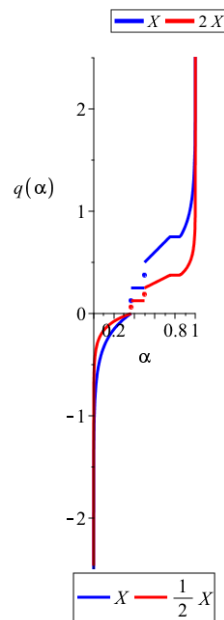
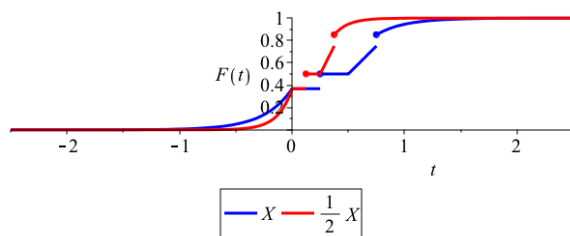
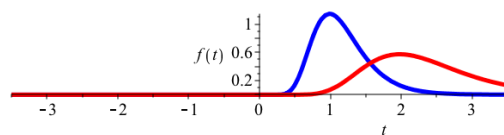
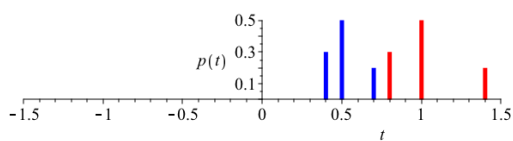
Vynásobení nenulovou konstantou r odpovídá podobnost ve směru vodorovné osy.

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

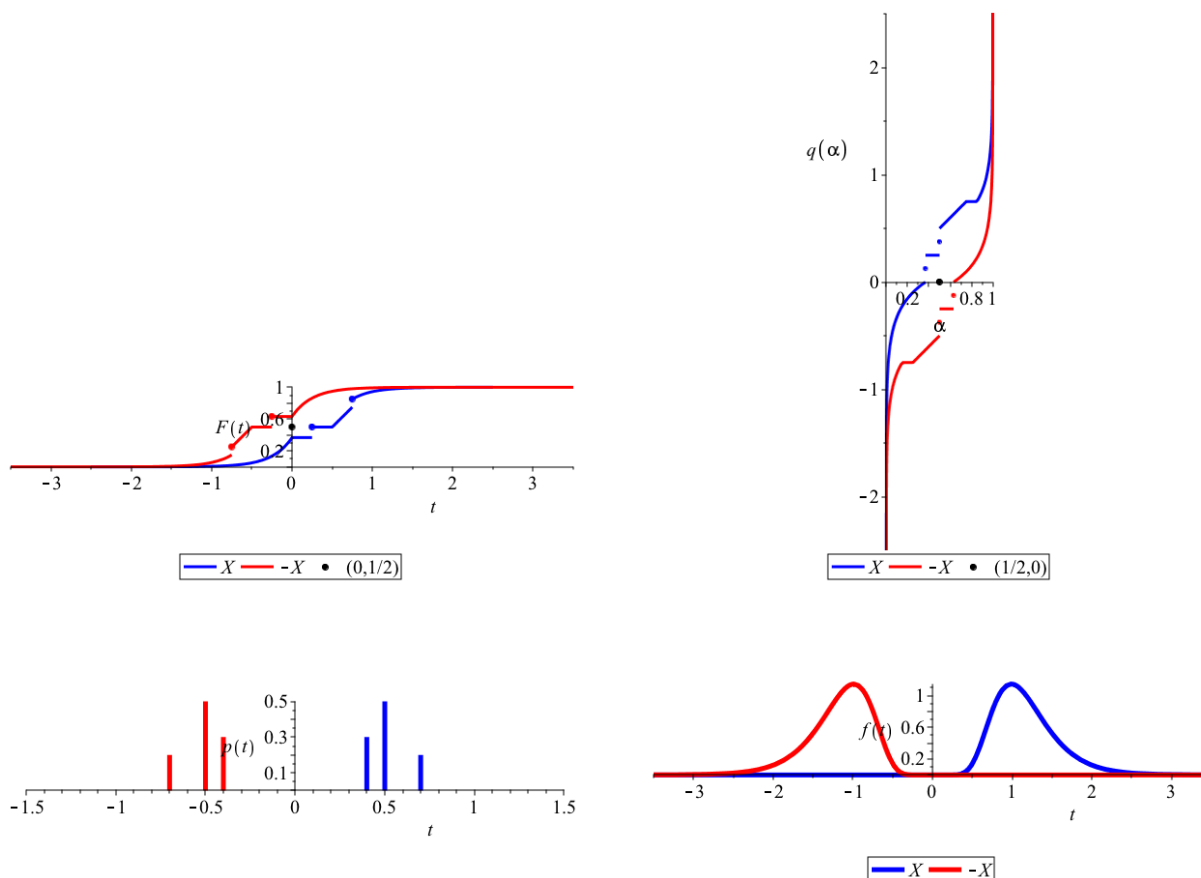
Pro distribuční funkci musíme rozlišit případy:

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





- $r = -1$: $F_{-X}(-t) = P(-X \leq -t) = P(X \geq t) = 1 - P(X < t)$, **v bodech spojitosti** distribuční funkce
 $F_{-X}(-t) = 1 - P(X < t) = 1 - P(X \leq t) = 1 - F_X(t)$,
 $F_{-X}(u) = 1 - F_X(-u)$, **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 :

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

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

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

Zobrazení nerostoucí funkcí h lze řešit jako zobrazení náhodné veličiny $-X$ neklesající funkcí $g(t) = h(-t)$.

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 rozděleními vstupních náhodných veličin a koeficienty směsi,

$$h(\text{Mix}_c(U, V)) = \text{Mix}_c(h(U), h(V)).$$

(Je jedno, jestli jakoukoli funkci h aplikujeme před, nebo po vytvoření směsi.)

4.8 Jak realizovat náhodnou veličinu na počítači

1. Vytvoříme náhodný (nebo pseudonáhodný) generátor náhodné veličiny X s rovnoměrným rozdělením na $\langle 0, 1 \rangle$. („Kolo štěstí.“)
2. Náhodná veličina $q_Y(X) \sim Y$ (kde \sim znamená „má stejné rozdělení“). Stačí na každou realizaci náhodné veličiny X aplikovat funkci q_Y . („Na kolo štěstí napíšeme vzestupně výsledky, každý s šířkou odpovídající jeho pravděpodobnosti.“)

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

5 Charakteristiky náhodných veličin

5.1 Střední hodnota

Motivační příklad (omezení kouření):

Chceme vyhodnotit, zda zákaz kouření v restauracích vedl k celkovému omezení kouření.

Místo mnoha hodnot chceme jednu „průměrnou“, která charakterizuje celou populaci.

Značení: E . nebo μ .

$$EX = \int_0^1 q_X(\alpha) d\alpha, \quad (\text{pokud existuje}).$$

Interpretace:

- Pro Diracovo rozdělení integrujeme konstantu přes interval délky 1.
- Pro diskrétní náhodnou veličinu U integrujeme po částech konstantní funkci; kde délka každého úseku je pravděpodobnost příslušné hodnoty,

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

- Obecný případ je limitou předchozího, pokud „diskretizační chyba“ $\rightarrow 0$.

Střední hodnota může být definována zvlášť pro

- **diskrétní** náhodnou veličinu U : $EU = \sum_{t \in \mathbb{R}} t \cdot p_U(t),$
- **spojitou** náhodnou veličinu V : $EV = \int_{-\infty}^{\infty} t \cdot f_V(t) dt,$
- **směs** náhodných veličin $X = \text{Mix}_c(V, U)$: $EX = cEV + (1-c)EU.$
(Může být V diskrétní, U spojitá. To **není** linearita střední hodnoty!)

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

Vzorec používající kvantilovou funkci lze zobecnit na střední hodnotu jakékoli (měřitelné) funkce h náhodné veličiny:

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

Speciálně pro **diskrétní** náhodnou veličinu U

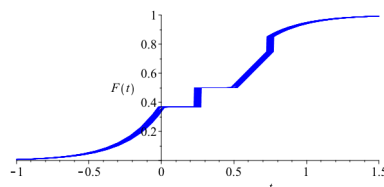
$$E(h(U)) = \sum_{t \in \Omega_U} h(t) \cdot p_U(t),$$

pro **spojitou** náhodnou veličinu V

$$E(h(V)) = \int_{-\infty}^{\infty} h(t) \cdot f_V(t) dt.$$

(Ale funkce spojitě náhodné veličiny nemusí být spojitá náhodná veličina.)

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 a je konečná** (což není vždy splněno).

Vlastnosti střední hodnoty

$$\begin{aligned} E r &= r, & \text{speciálně} & E(EX) = EX, \\ E(X + Y) &= EX + EY, & \text{speciálně} & E(X + r) = EX + r, \\ E(X - Y) &= EX - EY, \\ E(r X) &= r EX, & \text{obecněji} & E(r X + s Y) = r EX + s EY. \end{aligned}$$

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

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

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

Pouze pro **nezávislé** náhodné veličiny

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

kde \cdot značí součin nezávislých náhodných veličin.

5.2 Rozptyl (disperze)

Motivační příklad (stejně podnebí):

Chceme najít místo s podobným podnebím.

Průměrná teplota nestačí. Důležité je i kolísání teplot.

Značení: σ^2 , D ., var .

$$\begin{aligned} DX &= E((X - EX)^2) = E(X^2) - (EX)^2, \\ E(X^2) &= (EX)^2 + DX. \end{aligned} \tag{1}$$

Vlastnosti:

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

$$DX \geq 0,$$

$$D r = 0,$$

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

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

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

Pouze pro **nezávislé** náhodné veličiny ($\tilde{+}$, $\tilde{-}$ značí součet a rozdíl nezávislých)

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

5.3 Směrodatná odchylka

Značení: σ .

Má **stejný fyzikální rozměr** jako původní náhodná veličina (rozptyl nikoli).

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

Vlastnosti:

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

$$\sigma_X \geq 0,$$

$$\sigma_r = 0,$$

$$\sigma_{X+r} = \sigma_X,$$

$$\sigma_{rX} = |r| \sigma_X.$$

Pouze pro **nezávislé** náhodné veličiny

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

5.4 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:

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

$$E((X - EX)^k) = \int_0^1 (q_X(\alpha) - EX)^k d\alpha.$$

5.5 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

$$X = EX + \sigma_X \text{ norm } X.$$

Motivační příklad (biochemická vyšetření):

Laboratorní výsledky vydají mnoho čísel; abychom poznali, která jsou obvyklá a která znepokojivá, museli bychom znát alespoň jejich střední hodnoty a směrodatné odchylky.

Po znormování hned vidíme, které údaje zasluhují pozornost, aniž bychom museli studovat jejich typické hodnoty.

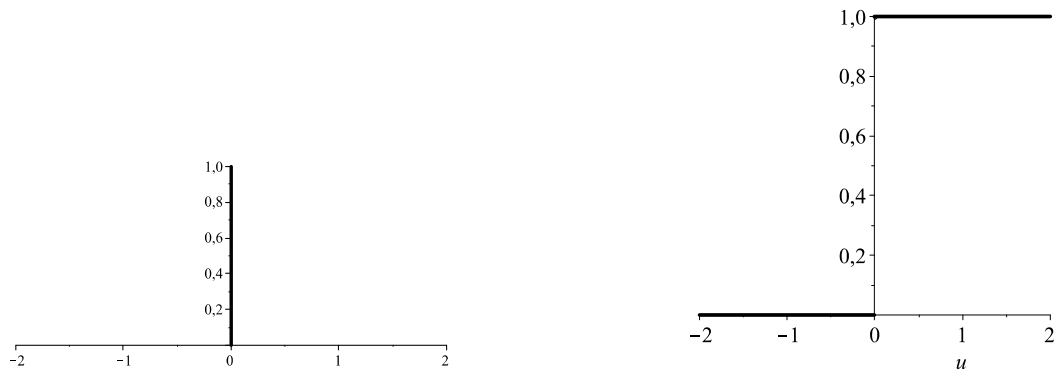
5.6 Základní typy diskrétních rozdělání

5.6.1 Diracova

s parametrem $r \in \mathbb{R}$, kterým je jediný možný výsledek,

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

Všechna diskrétní rozdělání jsou směsí Diracových rozdělání.



Pravděpodobnostní a distribuční funkce Diracova rozdělání pro $r = 0$

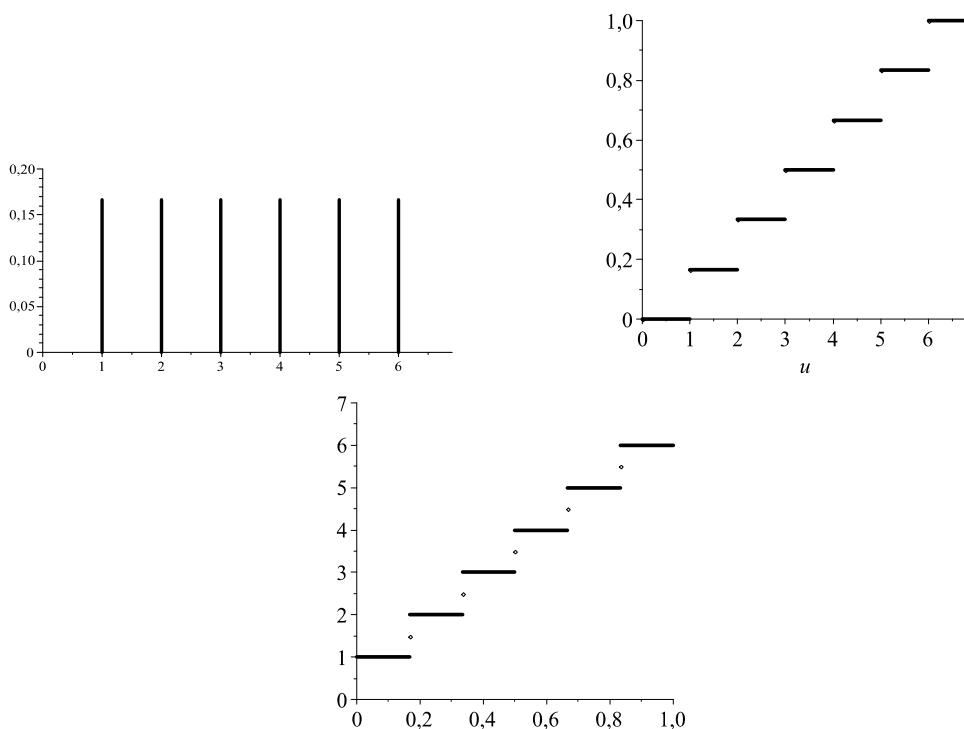
5.6.2 Rovnoměrná

s parametrem $m \in \mathbb{N}$, kterým je počet možných výsledků se stejnou pravděpodobností $1/m$.

Speciálně pro obor hodnot $\{1, 2, \dots, m\}$ dostáváme

$$p_X(k) = \frac{1}{m}, \quad k \in \{1, 2, \dots, m\},$$

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

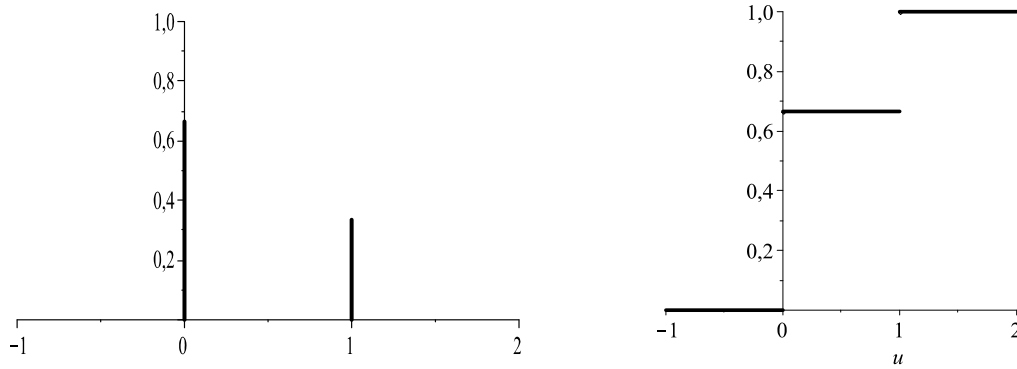


Pravděpodobnostní a distribuční funkce rovnoměrného rozdělání na $\{1, 2, \dots, 6\}$

5.6.3 Alternativní (Bernoulli) $\text{Alt}(q)$

s parametrem $q \in (0, 1)$, kterým je pravděpodobnost výsledku 1, s pravděpodobností $1 - q$ je výsledek 0, $\text{Alt}(q) = \text{Mix}_q(1, 0)$,

$$\begin{aligned} p_X(1) &= q, & p_X(0) &= 1 - q, \\ EX &= q, & DX &= q(1 - q). \end{aligned}$$



Pravděpodobnostní a distribuční funkce alternativního rozdělení pro $q = 1/3$

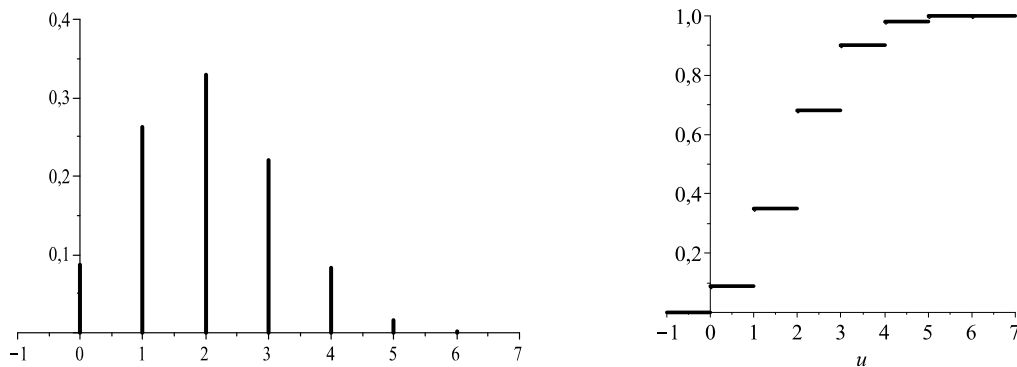
5.6.4 Binomická $\text{Bi}(m, q)$

s parametry $m \in \mathbb{N}$, $q \in (0, 1)$, jsou součtem m nezávislých náhodných veličin s alternativními rozděleními $\text{Alt}(q)$, $\text{Bi}(m, q) = \underbrace{\text{Alt}(q) + \dots + \text{Alt}(q)}_{m \times}$,

$$\text{Bi}(1, q) = \text{Alt}(q),$$

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

Součet **nezávislých** náhodných veličin s rozděleními $\text{Bi}(m, q)$, $\text{Bi}(n, q)$ má binomické rozdělení $\text{Bi}(m + n, q)$. Parametry by mohly být též m a $\lambda = m q = EX$, značíme $\text{Bi}^*(m, \lambda) = \text{Bi}(m, q)$. Výpočetní složitost výpočtu $p_X(k)$ je $O(k)$, celého rozdělení $O(m^2)$.



Pravděpodobnostní a distribuční funkce binomického rozdělení $\text{Bi}(6, 1/3)$

5.6.5 Poissonova $\text{Po}(\lambda)$

s parametrem $\lambda \in (0, \infty)$ jsou limitní případy binomických rozdělení $\text{Bi}(m, q_m) = \text{Bi}^*(m, \lambda)$ pro $m \rightarrow \infty$ při konstantním $\lambda = m q_m > 0$ (tedy $q_m \rightarrow 0$),

$$\text{Po}(\lambda) = \lim_{m \rightarrow \infty} \text{Bi}^*(m, \lambda).$$

Pravděpodobnostní funkce Poissonova rozdělení a binomických rozdělení se stejnou střední hodnotou 3:

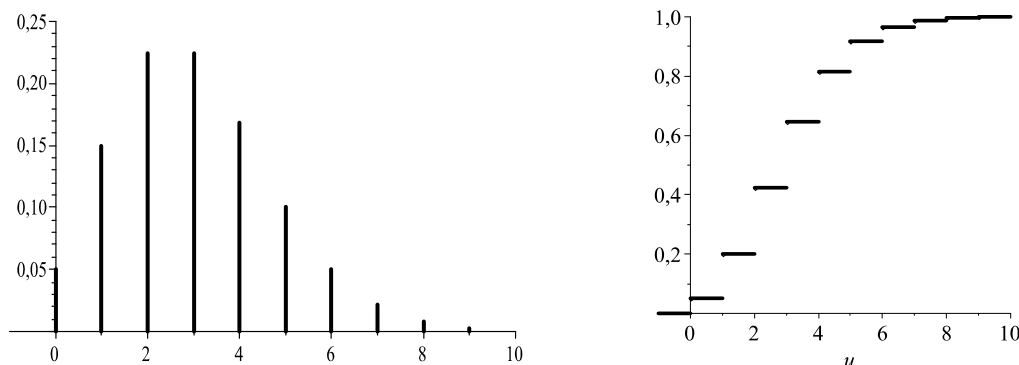
hodnota	0	1	2	3	4	5	6
$\text{Bi}^*(30, 3) = \text{Bi}(30, 0.1)$	0.042	0.141	0.228	0.236	0.177	0.102	0.047
$\text{Bi}^*(100, 3) = \text{Bi}(100, 0.03)$	0.047	0.147	0.225	0.227	0.171	0.101	0.050
$\text{Po}(3)$	0.050	0.149	0.224	0.224	0.168	0.101	0.050

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

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).



Pravděpodobnostní a distribuční funkce Poissonova rozdělení $\text{Po}(3)$

Poissonova rozdělení jako limitní případy binomických

Pro $m \rightarrow \infty$ při konstantním $m q_m = \lambda$, tj. $q_m = \frac{\lambda}{m}$:

$$\begin{aligned} p_{\text{Bi}(m, q_m)}(k) &= \binom{m}{k} q_m^k (1 - q_m)^{m-k} = \\ &= p_{\text{Bi}^*(m, \lambda)}(k) = \binom{m}{k} \left(\frac{\lambda}{m}\right)^k \left(1 - \frac{\lambda}{m}\right)^{m-k} \\ &\rightarrow \frac{\lambda^k}{k!} e^{-\lambda} = p_{\text{Po}(\lambda)}(k). \end{aligned}$$

Součet **nezávislých** náhodných veličin s Poissonovými rozděleními $\text{Po}(\lambda)$, $\text{Po}(\mu)$ má Poissonovo rozdělení $\text{Po}(\lambda + \mu)$.

5.6.6 Geometrická

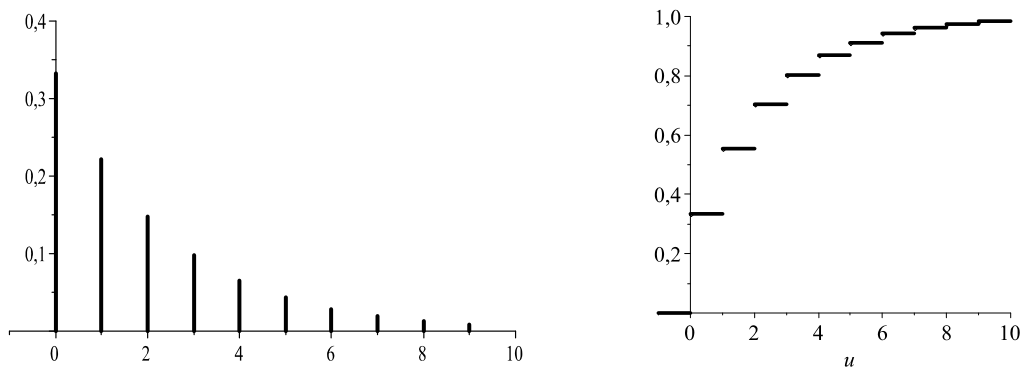
s parametrem $q \in (0, \infty)$ udávají

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

počet **neúspěchů** do prvního **úspěchu**, je-li v každém pokusu stejná pravděpodobnost **neúspěchu** $q \in (0, 1)$, neboli

počet jedniček na začátku posloupnosti nezávislých náhodných veličin s alternativními rozděleními $\text{Alt}(q)$,

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



Pravděpodobnostní a distribuční funkce geometrického rozdělení pro $q = 2/3$

5.6.7 Hypergeometrická

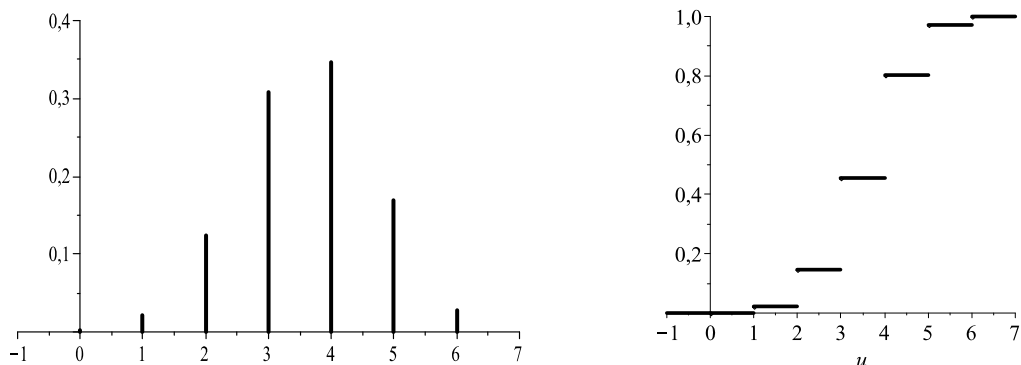
s parametry $M, J, m \in \mathbb{N}$, $m \leq J \leq M$.

Např. M losů, z nichž J vyhrává; náhodná veličina je počet vyhrávajících z m vybraných,

$$p_X(j) = \frac{\binom{J}{j} \binom{M-J}{m-j}}{\binom{M}{m}}, \quad j \in \{0, 1, 2, \dots, m\},$$

$$EX = \frac{mJ}{M}, \quad DX = \frac{mJ(M-J)(M-m)}{M^2(M-1)}.$$

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



Pravděpodobnostní a distribuční funkce hypergeometrického rozdělení pro $M = 25$, $J = 15$, $m = 6$

Binomická rozdělení jako limitní případy hypergeometrických

Pro $M \gg m$ je $\binom{M}{m} \doteq \frac{M^m}{m!}$ (Věta 2.5.6).

Hypergeometrická rozdělení pro $M \rightarrow \infty$ při konstantním $\frac{J_M}{M} = q$, tj. $\frac{M-J_M}{M} = 1-q$:

$$p_X(j) = \frac{\binom{J_M}{j} \binom{M-J_M}{m-j}}{\binom{M}{m}} \rightarrow \frac{\frac{J_M^j}{j!} \cdot \frac{(M-J_M)^{m-j}}{(m-j)!}}{\frac{M^m}{m!}} =$$

$$= \frac{m!}{j!(m-j)!} \cdot \frac{J_M^j}{M^j} \cdot \frac{(M-J_M)^{m-j}}{M^{m-j}} = \binom{m}{j} q^j (1-q)^{m-j} = p_{\text{Bi}(m,q)}(j).$$

5.7 Základní typy spojitých rozdělení

5.7.1 Rovnoměrná $R(a, b)$

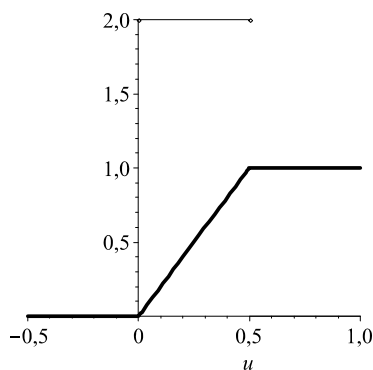
s parametry $a, b \in \mathbb{R}$, $a < b$.

$$f_X(t) = \begin{cases} \frac{1}{b-a} & \text{pro } t \in \langle a, b \rangle, \\ 0 & \text{jinak,} \end{cases}$$

$$F_X(u) = \begin{cases} \frac{u-a}{b-a} & \text{pro } u \in \langle a, b \rangle, \\ 0 & \text{pro } u < a, \\ 1 & \text{pro } u > b, \end{cases}$$

$$q_X(\alpha) = a + (b - a)\alpha,$$

$$EX = \frac{a+b}{2}, \quad DX = \frac{1}{12}(b-a)^2, \quad \sigma_X = \frac{1}{2\sqrt{3}}(b-a).$$



Hustota a distribuční funkce rovnoměrného rozdělení $R(0, 1/2)$

5.7.2 Normální (Gaussova) $N(\mu, \sigma^2)$

A. Normované $N(0, 1)$ s hustotou

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

Distribuční funkce je transcendentní (Gaussův integrál) Φ ,

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

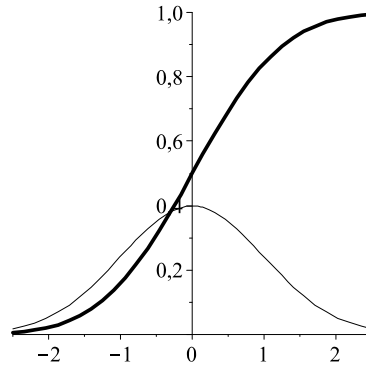
kvantilová funkce Φ^{-1} je inverzní k Φ .

B. Obecná $N(\mu, \sigma^2) = \mu + \sigma N(0, 1)$

s parametry $\mu \in \mathbb{R}$ (střední hodnota), $\sigma^2 \in (0, \infty)$ (rozptyl),

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

Součet dvou **nezávislých** veličin s normálním rozdělením $N(\mu_1, \sigma_1^2), N(\mu_2, \sigma_2^2)$ má normální rozdělení $N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$.



Hustota a distribuční funkce normovaného normálního rozdělení $N(0, 1)$

5.7.3 Logaritmicko-normální $LN(\mu, \sigma^2) = \exp(N(\mu, \sigma^2))$

s parametry $\mu \in \mathbb{R}$, $\sigma^2 \in (0, \infty)$ jsou rozdělení náhodných veličin $X = \exp(Y)$, kde Y má $N(\mu, \sigma^2)$,

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

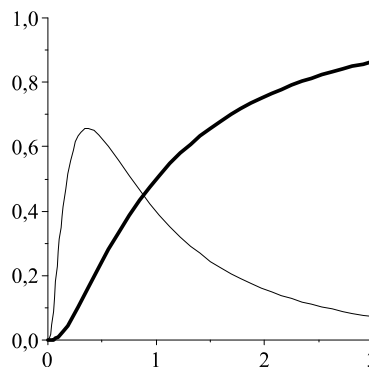
$$F_X(u) = \begin{cases} F_{N(\mu, \sigma^2)}(\ln u) & \text{pro } u > 0, \\ 0 & \text{jinak,} \end{cases}$$

$$EX = \exp\left(\mu + \frac{\sigma^2}{2}\right),$$

$$DX = (\exp(2\mu + \sigma^2)) (\exp(\sigma^2) - 1),$$

$$\sigma_X = \left(\exp\left(\mu + \frac{\sigma^2}{2}\right)\right) \sqrt{\exp(\sigma^2) - 1}.$$

Součin dvou **nezávislých** veličin s logaritmicko-normálním rozdělením $LN(\mu_1, \sigma_1^2)$, $LN(\mu_2, \sigma_2^2)$ má logaritmic-konormální rozdělení $LN(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$.

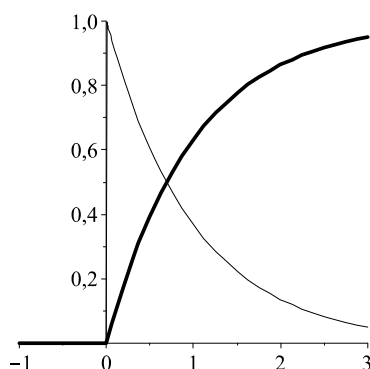


Hustota a distribuční funkce logaritmic-konormálního rozdělení $LN(0, 1)$

5.7.4 Exponenciální $Ex(\tau)$

s parametrem $\tau \in (0, \infty)$ (střední hodnota); 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 :

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



Hustota a distribuční funkce exponenciálního rozdělení $\text{Ex}(1)$

5.8 Čebyševova nerovnost

Motivační příklad (10 000 hodů mincí):

Při 10 000 hodech mincí má počet líců X binomické rozdělení $\text{Bi}(10\,000, \frac{1}{2})$ se střední hodnotou $EX = 5\,000$,

směrodatnou odchylkou $\sigma_X = \sqrt{10\,000 \cdot \frac{1}{4}} = 50$.

Jaká je pravděpodobnost, že se výsledek bude lišit od střední hodnoty o nejméně $200 = 4\sigma_X$?

$$\begin{aligned}
&\sum_{k=0}^{4\,800} p_{\text{Bi}(10\,000, 0.5)}(k) + \sum_{k=5\,200}^{10\,000} p_{\text{Bi}(10\,000, 0.5)}(k) = \\
&= 1 - \sum_{k=4\,801}^{5\,199} p_{\text{Bi}(10\,000, 0.5)}(k) = 1 - \sum_{k=4\,801}^{5\,199} \binom{10\,000}{k} \frac{1}{2^{10\,000}} \doteq 0.0000659.
\end{aligned}$$

Věta:

$$\forall \delta > 0 : P(|\text{norm } X| \geq \delta) \leq \frac{1}{\delta^2},$$

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

Důkaz pomocí kvantilové funkce:

$$\begin{aligned}
1 &= D(\text{norm } X) = E((\text{norm } X)^2) - \underbrace{(E(\text{norm } X))^2}_0 = \int_0^1 (q_{\text{norm } X}(\alpha))^2 d\alpha \\
&\geq \int_I (q_{\text{norm } X}(\alpha))^2 d\alpha,
\end{aligned}$$

kde $I = \{\alpha \in (0, 1) : |q_{\text{norm } X}(\alpha)| \geq \delta\}$ jsou 2 intervaly o celkové délce $P(|\text{norm } X| \geq \delta)$,

$$1 \geq \int_I (q_{\text{norm } X}(\alpha))^2 d\alpha \geq \int_I \delta^2 d\alpha = \delta^2 P(|\text{norm } X| \geq \delta).$$

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

$$\begin{aligned} \forall \delta > 0 : P(|\text{norm } X| < \delta) &\geq 1 - \frac{1}{\delta^2}, \\ \forall \delta > 0 : P\left(\left|\frac{X - EX}{\sigma_X}\right| \geq \delta\right) &\leq \frac{1}{\delta^2}, \\ \forall \varepsilon > 0 : P(|X - EX| \geq \varepsilon) &\leq \frac{\sigma_X^2}{\varepsilon^2} = \frac{DX}{\varepsilon^2}, \\ \forall \varepsilon > 0 : P(|X - EX| < \varepsilon) &\geq 1 - \frac{\sigma_X^2}{\varepsilon^2} = 1 - \frac{DX}{\varepsilon^2}. \end{aligned}$$

Motivační příklad (10 000 hodů mincí – pokračování):

Při 10 000 hodech mincí má počet líců X binomické rozdělení $\text{Bi}(10\,000, 0.5)$ se střední hodnotou $EX = 5\,000$,

směrodatnou odchylkou $\sigma_X = \sqrt{10\,000 \cdot 0.25} = 50$.

Jaká je pravděpodobnost odchylky od střední hodnoty o více než $4\sigma_X$?

$$\text{nejvýše } \left(\frac{1}{4}\right)^2 = \frac{1}{16}.$$

Buffonova úloha – přesnost odhadu

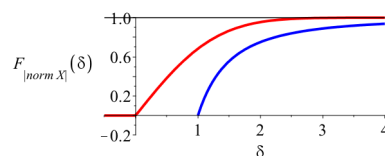
Při n hodech jehlou má počet protnutí linky X binomické rozdělení $\text{Bi}(n, \frac{2}{\pi})$.

Kolik hodů potřebujeme, abychom π odhadli na 99 % s přesností 1 %?

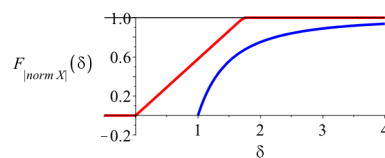
Uf! Šetřme elektrinu, ne hlavu!

$$\begin{aligned} 0.01 &\geq \frac{\sigma_X^2}{\varepsilon^2} = \frac{\sigma_X^2}{(0.01 EX)^2} = \frac{\frac{2}{\pi} (1 - \frac{2}{\pi}) n}{(0.01 \frac{2}{\pi} n)^2} = \frac{10\,000 (\frac{\pi}{2} - 1)}{n}, \\ n &\geq 1\,000\,000 \left(\frac{\pi}{2} - 1\right) \doteq 570\,796. \end{aligned}$$

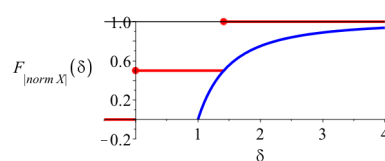
Uf! Šetřme ruce, ne hlavu!



Distribuční funkce absolutní hodnoty normovaného normálního rozdělení $|N(0, 1)|$ (červeně) ve srovnání s mezí dle Čebyševovy nerovnosti (modře)



Distribuční funkce absolutní hodnoty normovaného spojitého rovnoměrného rozdělení $|R(-\sqrt{3}, \sqrt{3})|$ (červeně) ve srovnání s mezí dle Čebyševovy nerovnosti (modře)



Distribuční funkce absolutní hodnoty normovaného binomického rozdělení $|\text{norm Bi}(2, 0.5)|$ (červeně) ve srovnání s mezí dle Čebyševovy nerovnosti (modře)

6 Náhodné vektory

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 $I_k = (-\infty, t_k]$, $t_k \in \mathbb{R}$,

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

$F_{\mathbf{X}}: \mathbb{R}^n \rightarrow [0, 1]$ je **sdrúžená 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**.

Postačující podmínky bychom dostali, kdybychom např. pro dvě dimenze přidali podmínky

$$0 \leq P_{\mathbf{X}}((a, b] \times (c, d]) = F_{\mathbf{X}}(b, d) - F_{\mathbf{X}}(a, d) - F_{\mathbf{X}}(b, c) + F_{\mathbf{X}}(a, c),$$

pro všechna $a, b, c, d \in \mathbb{R}$, $a < b$, $c < d$.

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

Podmínky nezávislosti pro složky náhodného vektoru jsou

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

obecněji

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

6.1 Diskrétní náhodný vektor

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

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

která 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:

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

6.2 Spojitý náhodný vektor

má všechny složky spojité. Lze jej popsat též **sduženou 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

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

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

6.3 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.)

6.4 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: $X = \Re(X) + i\Im(X)$: $EX := E\Re(X) + iE\Im(X)$
- 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

$$E(aU + b) = aEU + b, \quad D(aU + b) = a^2 DU.$$

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.

$$\begin{aligned} E(X + Y) &= EX + EY, \\ D(X + Y) &= E((X + Y)^2) - (E(X + Y))^2 \\ &= E(X^2 + Y^2 + 2XY) - (EX + EY)^2 \\ &= E(X^2) + E(Y^2) + 2E(XY) - ((EX)^2 + (EY)^2 + 2EXEY) \\ &= \underbrace{E(X^2) - (EX)^2}_{DX} + \underbrace{E(Y^2) - (EY)^2}_{DY} + 2 \underbrace{(E(XY) - EXEY)}_{\text{cov}(X,Y)} \\ &= DX + DY + 2 \text{cov}(X, Y), \end{aligned}$$

kde $\text{cov}(X, Y) := E(XY) - EXEY$ je **kovariance** náhodných veličin X, Y , též

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

neboť

$$\begin{aligned} E((X - EX)(Y - EY)) &= E(XY - XEY - YEX + EXEY) \\ &= E(XY) - EXEY - \underbrace{EYEX}_{EXEY} + EXEY \\ &= E(XY) - EXEY. \end{aligned}$$

Pro existenci kovariance je postačující existence rozptylů DX, DY .

Vlastnosti kovariance:

$$\begin{aligned} \text{cov}(X, X) &= DX, \quad \text{cov}(Y, X) = \text{cov}(X, Y), \\ \text{cov}(aX + b, cY + d) &= ac \text{cov}(X, Y) \quad (a, b, c, d \in \mathbb{R}) \\ (\text{srovnajte s vlastnostmi rozptylu jako speciálního případu}), \\ \text{speciálně} \quad \text{cov}(X, -X) &= -DX. \end{aligned}$$

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

$$\begin{aligned} \varrho(X, X) &= 1, \quad \varrho(X, -X) = -1, \quad \varrho(X, Y) \in \langle -1, 1 \rangle, \\ \varrho(Y, X) &= \varrho(X, Y), \\ \varrho(aX + b, cY + d) &= \text{sign}(ac) \varrho(X, Y) \quad (a, b, c, d \in \mathbb{R}, \quad a \neq 0 \neq c) \\ (\text{až na znaménko nezáleží na prosté lineární transformaci}). \end{aligned}$$

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

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.

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

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

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

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

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

6.5 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).

7 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,

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

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

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 konstantních náhodných veličin (tj. s Diracovým rozdělením),

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

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

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

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

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

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

$$\begin{aligned} X \bullet X &= \|X\|^2 = \|X - EX\|^2 + \|EX\|^2, \\ E(X^2) &= DX + (EX)^2. \end{aligned}$$

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

Speciálně pro náhodné veličiny z \mathcal{N} :

$$\begin{aligned} \sigma_X^2 &= X \bullet X, \\ \sigma_X &= \|X\|, \\ \text{cov}(X, Y) &= X \bullet Y, \\ \varrho(X, Y) &= \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{X \bullet Y}{\|X\| \|Y\|} = \cos \angle(X, Y). \end{aligned}$$

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} ,

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

POZOR! Nepletejte **nezávislost náhodných veličin** s **lineární nezávislostí** v lineárním prostoru, který tvoří!

7.2 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**:

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

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

8 Základní pojmy statistiky

8.1 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].

8.2 Náhodný výběr, empirické rozdělení

Motivační úloha: Máme odhadnout rozdělení hmotnosti lidí a jeho parametry.

Náhodný pokus X: Ze **základního souboru (=populace)** náhodně vybereme jeden objekt (člověka), ω (s rovnoměrným rozdělením).

Uřídíme jeho hmotnost, $X(\omega)$, kterou prozatím považujeme za konstantní, tj. s Diracovým rozdělením a rovnou její realizaci, $x(\omega) \in \mathbb{R}$.

Za množinu všech elementárních jevů Ω můžeme vzít množinu objektů (lidí), z nichž vybíráme, za σ -algebru $\exp \Omega$, pravděpodobnostní míra je $M \mapsto \frac{|M|}{|\Omega|}$, $M \subseteq \Omega$.

Náhodná veličina $X: \Omega \rightarrow \mathbb{R}$ (závisí jen výběru člověka) má diskrétní rozdělení (směs Diracových),

$$\begin{aligned} EX &= \frac{1}{|\Omega|} \sum_{\omega \in \Omega} x(\omega), \\ DX &= \frac{1}{|\Omega|} \sum_{\omega \in \Omega} (x(\omega) - EX)^2. \end{aligned}$$

Náhodný pokus $n \times X$: Náhodný pokus X nezávisle opakujeme n -krát.

Vybereme (s opakováním a s rovnoměrným rozdělením) n -tici lidí $(\omega_1, \dots, \omega_n) \in \Omega^n$. Výsledkem je n -tice hmotností

$$(x_1, \dots, x_n) := (x(\omega_1), \dots, x(\omega_n)).$$

Za množinu všech elementárních jevů můžeme vzít

$$\Omega^n = \{(\omega_1, \dots, \omega_n) \mid \omega_1 \in \Omega, \dots, \omega_n \in \Omega\},$$

za σ -algebru $\exp(\Omega^n)$, pravděpodobnostní míra je $M \mapsto \frac{|M|}{|\Omega|^n}$, $M \subseteq \Omega^n$.

Při j -tém opakování pokusu realizujeme náhodnou veličinu $X_j: \Omega^n \rightarrow \mathbb{R}$,

$$X_j(\omega_1, \dots, \omega_n) \mapsto X(\omega_j) = x_j,$$

která závisí jen na j -té souřadnici a je konstantní na všech elementárních jevech (n -ticích lidí), které ji mají stejnou (ve kterých j -tý člověk je stejný).

Celý pokus charakterizuje náhodný vektor $\mathbf{X} := (X_1, \dots, X_n): \Omega^n \rightarrow \mathbb{R}^n$,

$$(X_1, \dots, X_n)(\omega_1, \dots, \omega_n) \mapsto (x_1, \dots, x_n).$$

Všechny jeho složky jsou **nezávislé** a mají **stejné rozdělení** jako původní náhodná veličina X .

Původní rozdělení náhodné veličiny X **neznáme**, neboť jsme ho **nezměřili** na celém základním souboru. Posloupností n realizací jsme dostali **výběrový soubor**, který „reprezentuje“ základní soubor.

Výběrový soubor definuje obdobným způsobem tzv. **empirické rozdělení** $\text{Emp}(\mathbf{x})$:

Náhodný pokus Y: Z **výběrového souboru** $\mathbf{x} := (x_1, \dots, x_n)$ náhodně vybereme jeden prvek, j -tý (s rovnoměrným rozdělením). Výsledkem je $x_j \in \mathbb{R}$.

Empirické rozdělení známe, neboť jsme ho **změřili**.

Za množinu všech elementárních jevů Ω_{Emp} můžeme vzít $\{1, \dots, n\}$, za σ -algebru $\exp \Omega_{\text{Emp}}$, pravděpodobnostní míra je $M \mapsto \frac{|M|}{n}$, $M \subseteq \Omega_{\text{Emp}}$,

$$\begin{aligned} \mathbb{E} \text{Emp}(\mathbf{x}) &= \frac{1}{n} \sum_{j=1}^n x_j =: \bar{\mathbf{x}}, \\ \mathbb{D} \text{Emp}(\mathbf{x}) &= \frac{1}{n} \sum_{j=1}^n (x_j - \mathbb{E} \text{Emp}(\mathbf{x}))^2 = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{\mathbf{x}})^2. \end{aligned}$$

Tyto hodnoty mohou sloužit k odhadu neznámých parametrů původního rozdělení na základním souboru. Vše funguje obdobně, pokud pozorovaná náhodná veličina X není pro každého jedince konstantní (závisí na dalších parametrech, „okolnostech“).

Stačí v modelu použít více elementárních jevů, kterými mohou být např. dvojice (člověk, okolnosti). Závislost na okolnostech může být i spojitá.

Empirické rozdělení je vždy diskrétní (směs Diracových, $\text{Mix}_{(1/n, \dots, 1/n)}(x_1, \dots, x_n)$), používá již změřené veličiny a jediný náhodný vliv je výběr prvku z posloupnosti.

Obvykle vylučujeme vícenásobný výběr stejného prvku (*výběr bez vracení*), čímž se dopouštíme chyby. Ta se obvykle zanedbává, neboť

1. pro velký rozsah základního souboru není podstatná,
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**. Náhodný pokus $n \times X$ je jednou z možností, jak tyto podmínky zajistit.

(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**.

funkce $f: D \rightarrow \mathbb{R}$	funkční hodnota $f(x) \in \mathbb{R}, \quad x \in D$
náhodná veličina $X: \Omega \rightarrow \mathbb{R}$	realizace náhodné veličiny $x := X(\omega) \in \mathbb{R}, \quad \omega \in \Omega$
náhodný vektor/výběr $\mathbf{X} = (X_1, \dots, X_n): \Omega^n \rightarrow \mathbb{R}^n$	realizace náhodného vektoru/výběru $\mathbf{x} = (x_1, \dots, x_n) := \mathbf{X}(\omega) \in \mathbb{R}^n, \quad \omega \in \Omega^n$

Statistika je (každá) měřitelná funkce G , definovaná na náhodném výběru libovolného (dostatečného) 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 pro **odhad parametrů rozdělení** (které nám zůstávají skryté).

8.3 Obecné vlastnosti odhadů

Značení:

ϑ ... jakákoli hodnota parametru (reálné číslo),

ϑ^* ... skutečná (správná) hodnota parametru (reálné číslo),

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

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

Žádoucí vlastnosti odhadů:

- $E\hat{\Theta}_n = \vartheta^*$, tj. $E(\hat{\Theta}_n - \vartheta^*) = 0$ **nestranný** (opak: **vychýlený**)
- $\lim_{n \rightarrow \infty} E\hat{\Theta}_n = \vartheta^*$, tj. $\lim_{n \rightarrow \infty} E(\hat{\Theta}_n - \vartheta^*) = 0$ **asymptoticky nestranný**
- **eficientní** = s malým rozptylem, což posuzujeme podle $E((\hat{\Theta}_n - \vartheta^*)^2) = D\hat{\Theta}_n + (E\hat{\Theta}_n - \vartheta^*)^2$, 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)
- $\lim_{n \rightarrow \infty} E(\hat{\Theta}_n - \vartheta^*) = 0$, $\lim_{n \rightarrow \infty} \sigma_{\hat{\Theta}_n} = 0$ **konzistentní**
- **robustní**, tj. odolný vůči šumu („i při zašuměných datech dostáváme dobrý výsledek“) – přesné kritérium chybí, ale je to velmi praktická vlastnost

8.4 Odhad střední hodnoty

pomocí střední hodnoty empirického rozdělení (aritmetického průměru realizace náhodného výběru):

$$\bar{x} := E \text{Emp}(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n x_j.$$

Když totéž provedeme s náhodnými veličinami výběru, dostaneme náhodnou veličinu

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

\bar{X} = **výběrový průměr**, \bar{x} = **realizace výběrového průměru**.

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

Věta:

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

pokud existují. (Zde $EX = EX_j$ atd.) Výběrový průměr minimalizuje kritérium nejmenších čtverců

$$\ell_2(c) = E(c - \text{Emp}(\mathbf{x}))^2 = \frac{1}{n} \sum_{i=1}^n (c - x_i)^2.$$

Důsledek: Výběrový průměr je nestranný konzistentní odhad střední hodnoty.

(Nezávisle na typu rozdělení.)

Čebyševova nerovnost pro \bar{X}_n dává

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

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 $n DX \rightarrow \infty$. **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

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

konvergují k normovanému normálnímu rozdělení v následujícím smyslu:

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

neboli

$$\forall t \in \mathbb{R} : \lim_{n \rightarrow \infty} |F_{Y_n}(t) - \Phi(t)| = 0,$$

Pokud má původní rozdělení 3. centrální moment, je konvergence dokonce **stejněměrná**, tj.

$$\lim_{n \rightarrow \infty} \sup_{t \in \mathbb{R}} |F_{Y_n}(t) - \Phi(t)| = 0,$$

Motivační příklad (10 000 hodů mincí):

Jak malá je pravděpodobnost, že se výsledek bude lišit od střední hodnoty o nejméně $200 = 4\sigma_X$?

S pravděpodobností $1 - \Phi(4) \doteq 1 - 0.9999683 = 0.0000317$ bude aspoň o $4\sigma_X$ větší, se stejnou pravděpodobností o $4\sigma_X$ menší, celkem

$$2(1 - \Phi(4)) \doteq 2 \cdot 0.0000317 = 0.0000633,$$

v dobré shodě s přesnějším (a pracným) výsledkem 0.0000659.

Kdybychom zohlednili kvantizační chybu a uvažovali toleranci ± 199.5 , dostali bychom ještě lepší odhad 0.0000661.

Buffonova úloha – přesnost odhadu

Binomické rozdělení $\text{Bi}(n, \frac{2}{\pi})$ při n hodech se blíží normálnímu se střední hodnotou $EX = \frac{2}{\pi}n \doteq 0.6366n$,

směrodatnou odchylkou $\sigma_X = \sqrt{\frac{2}{\pi}(1 - \frac{2}{\pi})n} \doteq 0.48\sqrt{n}$.

Kolik hodů n potřebujeme, abychom π odhadli na 99 % s přesností 1 %?

S pravděpodobností 99 % normovaná hodnota bude v absolutní hodnotě menší než kvantil $\Phi^{-1}(0.995) \doteq 2.576$. To je tolerance měřená ve směrodatných odchylkách.

Má být nejvýše 1 % střední hodnoty,

$$\begin{aligned} 2.576\sigma_X &\leq 0.01EX, \\ 2.576\sqrt{\frac{2}{\pi}(1 - \frac{2}{\pi})n} &\leq 0.01\frac{2}{\pi}n, \\ \sqrt{n} &\geq 2.576 \cdot 100 \sqrt{\frac{\pi}{2} - 1}, \\ n &\geq 2.576^2 \cdot 10\,000 \left(\frac{\pi}{2} - 1\right) \doteq 37\,877. \end{aligned}$$

Výsledný počet bude v mezích $24\,113 \pm 241$ s pravděpodobností

$$\sum_{k=23\,872}^{24\,354} p_{\text{Bi}(37\,877, \frac{2}{\pi})}(k) = \sum_{k=23\,872}^{24\,354} \binom{37\,877}{k} \left(\frac{2}{\pi}\right)^k \left(1 - \frac{2}{\pi}\right)^{37\,877-k} \doteq 0.00988.$$

8.5 Odhad k -tého obecného momentu EX^k

pomocí k -tého obecného momentu empirického rozdělení:

$$m_{\mathbf{x}^k} := E \text{Emp}(\mathbf{x})^k = \frac{1}{n} \sum_{j=1}^n x_j^k \quad (\text{realizace výběrového } k\text{-tého obecného momentu}).$$

Je realizací odhadu

$$M_{\mathbf{X}^k} := \frac{1}{n} \sum_{j=1}^n X_j^k \quad (\text{výběrový } k\text{-tý obecný moment}).$$

Alternativní značení: M_k, m_k .

Věta: $EM_{\mathbf{X}^k} = EX^k$.

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

Důkaz:

$$EM_{\mathbf{X}^k} = E\left(\frac{1}{n} \sum_{j=1}^n X_j^k\right) = \frac{1}{n} \sum_{j=1}^n EX_j^k = EX^k,$$

$$DM_{\mathbf{X}^k} = \frac{1}{n^2} n DX^k = \frac{1}{n} (E(X^k)^2 - (EX^k)^2) = \frac{1}{n} (EX^{2k} - (EX^k)^2) \rightarrow 0.$$

8.6 Odhad rozptylu

8.6.1 Odhad rozptylu při známé střední hodnotě

$$\hat{\theta} = \frac{1}{n} \sum_{j=1}^n (x_j - EX)^2.$$

Je realizací odhadu

$$\hat{\Theta} = \frac{1}{n} \sum_{j=1}^n (X_j - EX)^2.$$

$$E\hat{\Theta} = \frac{1}{n} \sum_{j=1}^n \underbrace{E(X_j - EX)^2}_{DX} = DX.$$

- nestranný,
- konzistentní, pokud existuje 4. centrální moment

Rozdělení odhadu rozptylu pro výběr z normálního rozdělení $N(\mu, \sigma^2)$

$$X_j - \mu \sim N(0, \sigma^2) = \sigma N(0, 1)$$

$$U_j := \frac{X_j - \mu}{\sigma} \sim N(0, 1)$$

$$\sum_{j=1}^n U_j^2 = \frac{1}{\sigma^2} \sum_{j=1}^n (X_j - \mu)^2 = \frac{n \hat{\Theta}}{\sigma^2} \sim$$

8.6.2 Rozdělení χ^2 s n stupni volnosti, $\chi^2(n)$

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

$$E\chi^2(n) = E \sum_{j=1}^n U_j^2 = \sum_{j=1}^n \underbrace{EU_j^2}_{DU_j=1} = n, \quad E \frac{\chi^2(n)}{n} = 1,$$

$$D\chi^2(n) = D \sum_{j=1}^n U_j^2 = \sum_{j=1}^n \underbrace{DU_j^2}_2 = 2n, \quad D \frac{\chi^2(n)}{n} = \frac{2}{n}.$$

Součet **nezávislých** náhodných veličin s rozděleními $\chi^2(k)$, $\chi^2(n)$ má rozdělení $\chi^2(k+n)$.

Hustota

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

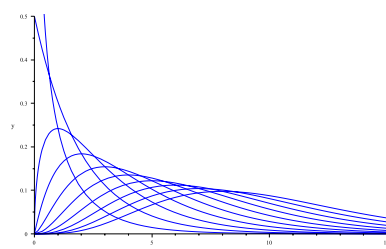
$$c(n) = \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})},$$

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

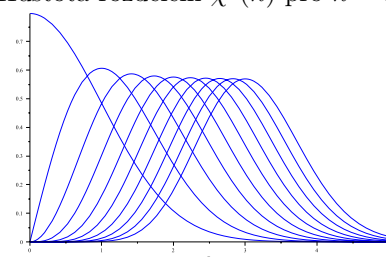
speciálně $\Gamma(m+1) = m!$ pro všechna $m \in \mathbb{N}$.

Speciálně pro $n = 2$ je $c(n) = 1/2$, $\chi^2(2) = \text{Ex}(2)$ (exponenciální rozdělení), $\frac{\chi^2(2)}{2} = \text{Ex}(1)$.

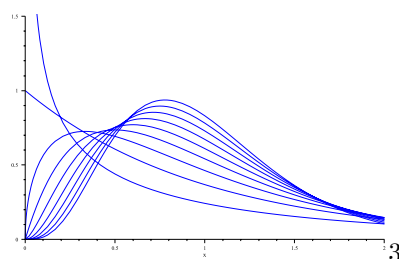
Důsledek: Součet m nezávislých náhodných veličin s exponenciálním rozdělením $\text{Ex}(2) = \chi^2(2)$ má rozdělení $\chi^2(2m)$.



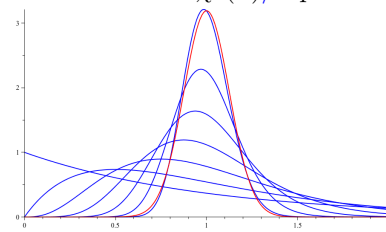
Hustota rozdělení $\chi^2(n)$ pro $n = 1, 2, \dots, 10$ stupňů volnosti.



Hustota rozdělení $\sqrt{\chi^2(n)}$ („vzdálenosti od středu terče“) pro $n = 1, 2, \dots, 10$ stupňů volnosti.



Hustota rozdělení $\chi^2(n)/n$ pro $n = 1, 2, \dots, 10$ stupňů volnosti.



Hustota rozdělení $\chi^2(n)/n$ pro $n = 2, 4, 8, \dots, 128$ stupňů volnosti.
Červeně náhrada posledního rozdělení normálním.

CLV pro $n \rightarrow \infty$

$$\Rightarrow \text{„}\chi^2(n) \rightsquigarrow N(n, 2n)\text{“}, \quad \text{„}\frac{\chi^2(n)}{n} \rightsquigarrow N(1, \frac{2}{n})\text{“},$$

přesněji $\text{norm } \chi^2(n) = \frac{\chi^2(n) - n}{\sqrt{2n}} \rightsquigarrow N(0, 1)$, $\sqrt{\frac{n}{2}} \left(\frac{\chi^2(n)}{n} - 1 \right) \rightsquigarrow N(0, 1)$.

Náhrada je „dobrá uprostřed“, kde je velká hustota.

Jenže my budeme potřebovat kvantily „na krajích“, a ty konvergují pomalu:

α -kvantily rozdělení $\chi^2(n)$ a jeho náhrady normálním rozdělením $N(n, 2n)$

$n \backslash \alpha$	0.005	0.01	0.025	0.05	0.95	0.975	0.99	0.995	0.999	0.9995
100	67.33	70.06	74.22	77.93	124.3	129.6	135.8	140.2	149.4	153.2
	63.57	67.10	72.28	76.74	123.3	127.7	132.9	136.4	143.7	146.5
200	152.2	156.4	162.7	168.3	234.0	241.1	249.4	255.3	267.5	272.4
	148.5	153.5	160.8	167.1	232.9	239.2	246.5	251.5	261.8	265.8
300	240.7	246.0	253.9	260.9	341.4	349.9	359.9	366.8	381.4	387.2
	236.9	243.0	252.0	259.7	340.3	348.0	357.0	363.1	375.7	380.6
400	330.9	337.2	346.5	354.6	447.6	457.3	468.7	476.6	493.1	499.7
	327.1	334.2	344.6	353.5	446.5	455.4	465.8	472.9	487.4	493.1

α -kvantily rozdělení $\chi^2(n)/n$ **vyděleného počtem stupňů volnosti** a jeho náhrady normálním rozdělením $N(1, 2/n)$

$n \backslash \alpha$	0.005	0.01	0.025	0.05	0.95	0.975	0.99	0.995	0.999	0.9995
100	0.673	0.701	0.742	0.779	1.24	1.30	1.36	1.40	1.49	1.53
	0.636	0.671	0.723	0.767	1.23	1.28	1.33	1.36	1.44	1.47
200	0.760	0.780	0.815	0.840	1.17	1.20	1.24	1.28	1.34	1.36
	0.742	0.767	0.804	0.836	1.16	1.20	1.23	1.26	1.31	1.33
300	0.803	0.820	0.847	0.870	1.14	1.17	1.20	1.22	1.27	1.29
	0.790	0.810	0.840	0.866	1.13	1.16	1.19	1.21	1.25	1.27
400	0.828	0.842	0.865	0.888	1.12	1.14	1.17	1.19	1.23	1.25
	0.818	0.836	0.861	0.884	1.12	1.14	1.16	1.18	1.22	1.23

8.6.3 Odhad rozptylu při **neznámé** střední hodnotě

pomocí rozptylu empirického rozdělení:

$$\widehat{\sigma_x^2} := D \text{Emp}(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2.$$

Je realizací odhadu

$$\widehat{\sigma_X^2} := \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X})^2 \neq \frac{1}{n} \sum_{j=1}^n (X_j - EX)^2.$$

Vzorce jsou dvouprůchodové, ale lze je upravit na jednorůchodové (numericky nevhodné):

$$\begin{aligned} \widehat{\sigma_X^2} &= \frac{1}{n} \left(\sum_{j=1}^n X_j^2 - 2 \bar{X} \underbrace{\sum_{j=1}^n X_j}_{n \bar{X}} + \underbrace{\sum_{j=1}^n \bar{X}^2}_{n \bar{X}^2} \right) = \frac{1}{n} \sum_{j=1}^n X_j^2 - \bar{X}^2, \\ \widehat{\sigma_x^2} &= \frac{1}{n} \sum_{j=1}^n x_j^2 - \bar{x}^2 = \frac{1}{n} \sum_{j=1}^n x_j^2 - \left(\frac{1}{n} \sum_{j=1}^n x_j \right)^2. \end{aligned}$$

Věta: $\widehat{\sigma_X^2}$ je **vychýlený** konzistentní odhad rozptylu (pokud původní rozdělení má rozptyl a 4. centrální moment).

Důkaz (pouze první části, s použitím jednorůchodového vzorce):

$$\begin{aligned} E \widehat{\sigma_X^2} &= E \left(\frac{1}{n} \sum_{j=1}^n X_j^2 - \bar{X}^2 \right) = EX^2 - E \bar{X}^2 = (EX)^2 + DX - \underbrace{(E \bar{X})^2}_{EX} - \underbrace{D \bar{X}}_{\frac{1}{n} DX} = \\ &= \left(1 - \frac{1}{n} \right) DX = \frac{n-1}{n} DX \rightarrow DX \quad \text{pro } n \rightarrow \infty. \end{aligned}$$

⇒ **nestranný** odhad:

$$S_{\mathbf{X}}^2 := \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{\mathbf{X}}_n)^2 = \frac{n}{n-1} \widehat{\sigma_X^2} \quad (\text{výběrový rozptyl}),$$

$$s_{\mathbf{x}}^2 := \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x}_n)^2 = \frac{n}{n-1} \widehat{\sigma_x^2} \quad (\text{realizace výběrového rozptylu}).$$

Alternativní značení: S^2, s^2 . (Dvojka v horním indexu neznamena kvadrát.) Jednoduchý vzorec – praktičtější, ale numericky horší:

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

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

Věta:

$$\text{ES}_{\mathbf{X}}^2 = \text{DX}.$$

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

Rozdělení výběrového rozptylu může být podstatně složitější.

Rozdělení výběrového rozptylu pro výběr z normálního rozdělení (dle [Likeš, Machek])

Problém: Ve výrazu

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

jsou veličiny $X_j - \bar{\mathbf{X}}$ závislé.

1. Pro $n = 2$ a X_1, X_2 s rozdělením $\text{N}(0, 1)$:

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

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

kde $U = \frac{X_1 - X_2}{\sqrt{2}}$ má rozdělení $\text{N}(0, 1)$, takže $S_{\mathbf{X}}^2$ má rozdělení $\chi^2(1)$.

2. Pro X_1, \dots, X_n s rozdělením $\text{N}(0, 1)$:

Náhodný vektor $\mathbf{X} = (X_1, \dots, X_n)$ s nezávislými složkami má rozdělení sféricky symetrické (kolem počátku), které se nezmění rotací ani ortonormální transformací souřadnic, $\mathbf{U} = \mathbf{X} \mathbf{M}$, kde $\mathbf{M} \in \mathbb{R}^{n \times n}$ je ortonormální matice. Použijeme ortogonální, resp. ortonormální matici

$$\begin{bmatrix} 1 & 1 & 1 & \cdots & 1 & 1 \\ -1 & 1 & 1 & \ddots & 1 & 1 \\ 0 & -2 & 1 & \ddots & 1 & 1 \\ 0 & 0 & -3 & \ddots & 1 & 1 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & -(n-1) & 1 \end{bmatrix}, \text{ resp. } \mathbf{M} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{12}} & \cdots & \frac{1}{\sqrt{n(n+1)}} & \frac{1}{\sqrt{n}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{12}} & \ddots & \frac{1}{\sqrt{n(n+1)}} & \frac{1}{\sqrt{n}} \\ 0 & -\frac{2}{\sqrt{6}} & \frac{1}{\sqrt{12}} & \ddots & \frac{1}{\sqrt{n(n+1)}} & \frac{1}{\sqrt{n}} \\ 0 & 0 & -\frac{3}{\sqrt{12}} & \ddots & \frac{1}{\sqrt{n(n+1)}} & \frac{1}{\sqrt{n}} \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & -\frac{n-1}{\sqrt{n(n+1)}} & \frac{1}{\sqrt{n}} \end{bmatrix}.$$

Postup pro $n = 2$ byl speciálním případem pro

$$\begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}, \text{ resp. } \mathbf{M} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}.$$

Transformací se zachovala nezávislost U_1, \dots, U_n , jejich rozdělení $N(0, 1)$ i součet

$$\sum_{j=1}^n X_j^2 = \|(X_1, \dots, X_n)\|^2 = \|(U_1, \dots, U_n)\|^2 = \sum_{j=1}^n U_j^2,$$

kde pro poslední člen platí $U_n = \sqrt{n} \bar{X}$, $U_n^2 = n \bar{X}^2$,

$$\sum_{j=1}^{n-1} U_j^2 = \sum_{j=1}^n X_j^2 - U_n^2 = \sum_{j=1}^n X_j^2 - n \bar{X}^2 = n \widehat{\sigma_X^2} = (n-1) S_X^2,$$

$$n \widehat{\sigma_X^2} = (n-1) S_X^2 \quad \text{má rozdělení} \quad \chi^2(n-1),$$

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

$$\widehat{\sigma_X^2} \quad \text{má rozdělení} \quad \frac{\chi^2(n-1)}{n}.$$

3. Pro X_1, \dots, X_n s rozdělením $N(\mu, \sigma^2)$:

μ nemá vliv, σ^2 se vše vynásobilo,

$$S_X^2 \sim \frac{\chi^2(n-1)}{n-1} \sigma^2,$$

$$\frac{S_X^2}{\sigma^2} \sim \frac{\chi^2(n-1)}{n-1},$$

$$\widehat{\sigma_X^2} \sim \frac{\chi^2(n-1)}{n} \sigma^2,$$

$$\frac{n \widehat{\sigma_X^2}}{\sigma^2} = \frac{(n-1) S_X^2}{\sigma^2} \sim \chi^2(n-1).$$

Důsledky:

$$E S_X^2 = \frac{n-1}{n-1} DX = DX \quad (\text{to už víme i obecně}),$$

$$D S_X^2 = \frac{2(n-1)}{(n-1)^2} (DX)^2 = \frac{2}{n-1} (DX)^2 \rightarrow 0 \quad \text{pro } n \rightarrow \infty.$$

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

8.6.4 Eficiency odhadů rozptylu pro **normální** rozdělení

1. efficiency odhadu S_X^2 (z vlastností rozdělení χ^2):

$$E(S_X^2 - DX)^2 = D S_X^2 = \frac{2}{n-1} (DX)^2.$$

2. efficiency odhadu $\widehat{\sigma_X^2}$ (DX je konstanta):

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

a protože

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

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

8.7 Odhad směrodatné odchylky

pomocí směrodatné odchylky empirického rozdělení:

$$\widehat{\sigma}_{\mathbf{x}} := \sigma_{\text{Emp}(\mathbf{x})} = \sqrt{\frac{1}{n} \sum_{j=1}^n (x_j - \bar{\mathbf{x}})^2}.$$

Je realizací odhadu

$$\widehat{\sigma}_{\mathbf{X}} := \sqrt{\frac{1}{n} \sum_{j=1}^n (X_j - \bar{\mathbf{X}})^2}.$$

Je vychýlený. Alternativa:

$$S_{\mathbf{X}} = \sqrt{S_{\mathbf{X}}^2} = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{\mathbf{X}}_n)^2} \quad (\text{výběrová směrodatná odchylka}),$$

$$s_{\mathbf{x}} = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{\mathbf{x}}_n)^2} \quad (\text{realizace výběrové směrodatné odchylky}).$$

Alternativní značení: S, s

Věta:

$$\mathbb{E} S_{\mathbf{X}} \leq \sigma_{\mathbf{X}}.$$

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

Důkaz:

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

$$\sigma_{\mathbf{X}} \geq \mathbb{E} S_{\mathbf{X}}.$$

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

8.8 Histogram a popis empirického rozdělení

V realizaci náhodného výběru $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$ nezáleží na pořadí složek (ale záleží na jejich opakování). Úsporněji je popsán množinou (nejvýše n) hodnot $H := \{x_1, \dots, x_n\}$ a jejich **četnostmi** n_t , $t \in H$. Data lze popsat **tabulkou četností** nebo grafem zvaným **histogram**.

Normováním dostaneme **relativní četnosti** $r_t := \frac{n_t}{n} = p_{\text{Emp}(\mathbf{x})}(t)$ (=hodnoty pravděpodobnostní funkce empirického rozdělení $\text{Emp}(\mathbf{x})$), $t \in H$, kde $\sum_{t \in H} r_t = 1$.

Výpočet z četností je jednodušší (pokud se opakují stejné hodnoty):

$$\begin{aligned} \mathbb{E} \text{Emp}(\mathbf{x}) &= \sum_{t \in H} t r_t = \frac{1}{n} \sum_{t \in H} t n_t = \frac{1}{n} \sum_{j=1}^n x_j = \bar{\mathbf{x}}, \\ \mathbb{E} (\text{Emp}(\mathbf{x}))^k &= \sum_{t \in H} t^k r_t = \frac{1}{n} \sum_{t \in H} t^k n_t = \frac{1}{n} \sum_{j=1}^n x_j^k = m_{\mathbf{x}^k}. \\ D \text{Emp}(\mathbf{x}) &= \sum_{t \in H} (t - \bar{\mathbf{x}})^2 r_t = \frac{1}{n} \sum_{t \in H} (t - \bar{\mathbf{x}})^2 n_t = \\ &= \frac{1}{n} \sum_{j=1}^n (x_j - \bar{\mathbf{x}})^2 = \widehat{\sigma}_{\mathbf{x}}^2 = \frac{n-1}{n} s_{\mathbf{x}}^2. \end{aligned}$$

8.9 Odhad mediánu

pomocí mediánu empirického rozdělení, $q_{\text{Emp}(\mathbf{x})}(\frac{1}{2})$ (**výběrový medián**). Poskytuje jinou informaci než výběrový průměr, mnohdy užitečnější (mj. **robustnější** – odolnější vůči vlivu vychýlených hodnot, *outliers*).

Výběrový medián minimalizuje kritérium

$$\ell_1(c) = \mathbb{E}|c - \text{Emp}(\mathbf{x})| = \frac{1}{n} \sum_{i=1}^n |c - x_i|.$$

Navíc víme, jak se změní monotónní funkcí h : $q_{\text{Emp}(h(\mathbf{x}))}(\frac{1}{2}) = h(q_{\text{Emp}(\mathbf{x})}(\frac{1}{2}))$.
Proč se používá méně než výběrový průměr:

- Vyšší výpočetní náročnost: seřazení hodnot má pracnost úměrnou $n \ln n$, výběrový průměr n , a i když existují algoritmy počítající medián s lineární složitostí, jsou komplikované.
- Vyšší paměťová náročnost: úměrná n , u výběrového průměru stačí 2 registry.
- Obtížná decentralizace a paralelizace výpočtu.

Obecněji lze odhadnout α -kvantil $q_X(\alpha)$ pomocí α -kvantilu empirického rozdělení, $q_{\text{Emp}(\mathbf{x})}(\alpha)$. Nesmíme však volit α blízké 0 nebo 1, nemůžeme např. na základě výběru rozsahu 1000 odhadovat kvantil $q_X(10^{-6})$. Pokud známe model rozdělení, můžeme jeho parametry odhadnout pomocí **všech** hodnot ve výběrovém souboru a z parametrů pak odhadnout požadované kvantily.

8.10 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 $\alpha \in (0, 1)$ je pravděpodobnost, že meze intervalu I budou překročeny; $1 - \alpha$ je **koefficient spolehlivosti**. Obvykle hledáme **horní**, resp. **dolní jednostranný** odhad,

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

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

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

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

8.11 Intervalové odhady parametrů **normálního** rozdělení

Předpoklad: Náhodná veličina $X \sim N(\mu, \sigma^2)$.

8.11.1 Intervalový odhad střední hodnoty při **známém** rozptylu σ^2

σ^2 známe; X odhadneme z realizace x ; μ neznáme

S pravděpodobností $1 - \alpha$ je $q_{N(\mu, \sigma^2)}(\frac{\alpha}{2}) \leq X \leq q_{N(\mu, \sigma^2)}(1 - \frac{\alpha}{2})$,
po znormování

$$-\Phi^{-1}(1 - \frac{\alpha}{2}) = \Phi^{-1}(\frac{\alpha}{2}) \leq \text{norm } X = \frac{X - \mu}{\sigma} \leq \Phi^{-1}(1 - \frac{\alpha}{2}),$$
$$X - \sigma \Phi^{-1}(1 - \frac{\alpha}{2}) \leq \mu \leq X + \sigma \Phi^{-1}(1 - \frac{\alpha}{2}).$$

Obvykle místo jedné realizace použijeme realizaci výběrového průměru $\bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right)$; po normalizaci

$$\text{norm } \bar{X}_n = \frac{\sqrt{n}}{\sigma} (\bar{X}_n - \mu) \sim N(0, 1).$$

S pravděpodobností $1 - \alpha$ je

$$\bar{X}_n - \frac{\sigma}{\sqrt{n}} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \leq \mu \leq \bar{X}_n + \frac{\sigma}{\sqrt{n}} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right),$$

pro jednostranné odhady

$$\mu \leq \bar{X}_n + \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1 - \alpha),$$

$$\bar{X}_n - \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1 - \alpha) \leq \mu,$$

Dostali jsme intervalové odhady pro μ

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

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

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má (nenulový) rozptyl a rozsah výběru je velký.

8.11.2 Intervalový odhad střední hodnoty při neznámém rozptylu

σ^2 neznáme, ale odhadneme z realizace výběrového rozptylu $S_X^2 \sim \frac{\sigma^2}{n-1} \chi^2(n-1)$;

$\bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ odhadneme z realizace \bar{x}_n ;

μ neznáme

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

8.11.3 Studentovo t-rozdělení (autor: Gossett)

s n stupni volnosti je rozdělení náhodné veličiny $\frac{U}{\sqrt{\frac{V}{n}}}$ $\left(= \frac{U}{\sqrt{W}} \right)$,

kde $U \sim N(0, 1)$,

$V \sim \chi^2(n)$ $\left(\frac{V}{n} = W \sim \frac{\chi^2(n)}{n} \right)$,
 U, V jsou nezávislé $(U, W \text{ jsou nezávislé})$.

Značení: $t(n)$.

Hustota:

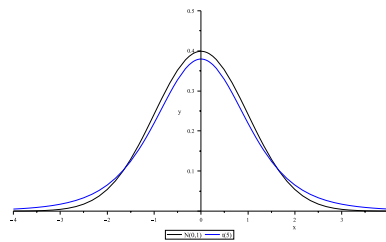
$$f_{t(n)}(x) = c(n) \left(1 + \frac{x^2}{n} \right)^{-\frac{1+n}{2}}, \quad c(n) = \frac{\Gamma(\frac{1+n}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})}.$$

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

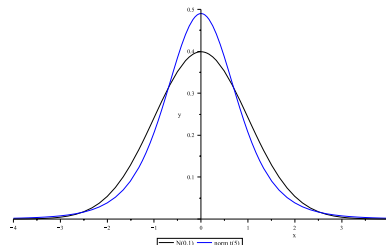
$t(1)$ je Cauchyho rozdělení, které nemá střední hodnotu,

$$f_{t(1)}(x) = \frac{1}{\pi} \frac{1}{1 + x^2}.$$

Pro velký počet stupňů volnosti se nahrazuje normovaným normálním rozdělením.



Hustota normovaného normálního rozdělení a Studentova rozdělení s 5 stupni volnosti.



Hustota normovaného normálního rozdělení a **normovaného** Studentova rozdělení s 5 stupni volnosti.

8.11.4 Intervalový odhad střední hodnoty při **neznámém** rozptylu 2

V našem případě:

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

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

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

Z toho vyplývají intervalové odhady

$$\left\langle \bar{X}_n - \frac{S_{\mathbf{X}}}{\sqrt{n}} q_{t(n-1)}\left(1 - \frac{\alpha}{2}\right), \bar{X}_n + \frac{S_{\mathbf{X}}}{\sqrt{n}} q_{t(n-1)}\left(1 - \frac{\alpha}{2}\right) \right\rangle,$$

$$\left\langle -\infty, \bar{X}_n + \frac{S_{\mathbf{X}}}{\sqrt{n}} q_{t(n-1)}(1 - \alpha) \right\rangle,$$

$$\left\langle \bar{X}_n - \frac{S_{\mathbf{X}}}{\sqrt{n}} q_{t(n-1)}(1 - \alpha), \infty \right\rangle.$$

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

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má nenulový rozptyl a rozsah výběru je velký (pak můžeme místo Studentova rozdělení použít normální).

8.11.5 Intervalový odhad rozptylu

μ nás nezajímá;

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

$\frac{(n-1) S_{\mathbf{X}}^2}{\sigma^2} \sim \chi^2(n-1)$.

S pravděpodobností $1 - \alpha$ je

$$q_{\chi^2(n-1)}\left(\frac{\alpha}{2}\right) \leq \frac{(n-1) S_{\mathbf{X}}^2}{\sigma^2} \leq q_{\chi^2(n-1)}\left(1 - \frac{\alpha}{2}\right),$$

$$\frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}\left(1 - \frac{\alpha}{2}\right)} \leq \sigma^2 \leq \frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}\left(\frac{\alpha}{2}\right)},$$

jednostranné odhady

$$\sigma^2 \leq \frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}(\alpha)},$$

$$\frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}(1-\alpha)} \leq \sigma^2.$$

Alternativní formulace:

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

$$\frac{S_{\mathbf{X}}^2}{\sigma^2} \sim \frac{\chi^2(n-1)}{n-1}.$$

S pravděpodobností $1 - \alpha$ je

$$q_{\frac{\chi^2(n-1)}{n-1}}\left(\frac{\alpha}{2}\right) \leq \frac{S_{\mathbf{X}}^2}{\sigma^2} \leq q_{\frac{\chi^2(n-1)}{n-1}}\left(1 - \frac{\alpha}{2}\right),$$

$$\frac{S_{\mathbf{X}}^2}{q_{\frac{\chi^2(n-1)}{n-1}}\left(1 - \frac{\alpha}{2}\right)} \leq \sigma^2 \leq \frac{S_{\mathbf{X}}^2}{q_{\frac{\chi^2(n-1)}{n-1}}\left(\frac{\alpha}{2}\right)},$$

jednostranné odhady

$$\sigma^2 \leq \frac{S_{\mathbf{X}}^2}{q_{\frac{\chi^2(n-1)}{n-1}}(\alpha)},$$

$$\frac{S_{\mathbf{X}}^2}{q_{\frac{\chi^2(n-1)}{n-1}}(1-\alpha)} \leq \sigma^2.$$

Dostali jsme intervalové odhady pro σ^2

$$\left\langle \frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}\left(1 - \frac{\alpha}{2}\right)}, \frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}\left(\frac{\alpha}{2}\right)} \right\rangle, \quad \left\langle \frac{S_{\mathbf{X}}^2}{q_{\frac{\chi^2(n-1)}{n-1}}\left(1 - \frac{\alpha}{2}\right)}, \frac{S_{\mathbf{X}}^2}{q_{\frac{\chi^2(n-1)}{n-1}}\left(\frac{\alpha}{2}\right)} \right\rangle,$$

$$\left(-\infty, \frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}(\alpha)}\right), \quad \left(-\infty, \frac{S_{\mathbf{X}}^2}{q_{\frac{\chi^2(n-1)}{n-1}}(\alpha)}\right),$$

$$\left(\frac{(n-1) S_{\mathbf{X}}^2}{q_{\chi^2(n-1)}(1-\alpha)}, \infty\right), \quad \left(\frac{S_{\mathbf{X}}^2}{q_{\frac{\chi^2(n-1)}{n-1}}(1-\alpha)}, \infty\right).$$

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

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má nenulový rozptyl a rozsah výběru je velký (pak můžeme místo rozdělení $\chi^2(n-1)$ použít normální $N(n-1, 2(n-1))$).

8.11.6 Intervalové odhady spojitých rozdělení, která nejsou normální

převádíme obvykle na normální rozdělení nelineární neklesající transformací

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

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

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

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

Někdy se transformuje na jiné rozdělení, např. Studentovo s vhodným počtem stupňů volnosti.

Poznámka: Je možné najít platné intervalové odhady, i když neexistuje střední hodnota nebo rozptyl.

8.12 Obecné odhady parametrů

Motivační příklad (volební předpověď):

Máme odhadnout výsledky voleb (kompletní).

Dosavadní postupy nám dovolovaly pouze odhad (včetně intervalového) výsledků jedné strany.

Hledanými parametry mohou být výsledky všech zúčastněných stran, vyjádřené vektorem čísel z $\langle 0, 1 \rangle$, jejichž součet je 1.

Motivační příklad (směs normálních rozdělání):

Spojité rozdělání, jehož hustota má více maxim, aproximujeme směsí normálních,

$$\text{Mix}_{(c_1, \dots, c_m)}(\mathcal{N}(\mu_1, \sigma_1^2), \dots, \mathcal{N}(\mu_m, \sigma_m^2)),$$

s hustotou

$$\sum_{i=1}^m c_i \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(-\frac{(t - \mu_i)^2}{2\sigma_i^2}\right).$$

Vektor parametrů je $(c_1, \dots, c_m; \mu_1, \dots, \mu_m; \sigma_1, \dots, \sigma_m)$ s omezujícími podmínkami

$$\begin{aligned} \sigma_i &\geq 0, & i &= 1, \dots, m, \\ 0 \leq c_i &\leq 1, & i &= 1, \dots, m, \\ \sum_{i=1}^m c_i &= 1. \end{aligned}$$

Formulace úlohy

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 $\mathbf{x} = (x_1, \dots, x_n)$.

8.12.1 Metoda momentů

(angl. *moment matching*)

Pro $k = 1, 2, \dots$ lze k -tý obecný moment vypočítat z modelu jako funkci $\boldsymbol{\vartheta}$,

$$\mathbb{E}X^k(\boldsymbol{\vartheta}) = \mathbb{E}X^k(\vartheta_1, \dots, \vartheta_i)$$

a současně odhadnout pomocí výběrového k -tého obecného momentu

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

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

$$\mathbb{E}X^k(\hat{\vartheta}_1, \dots, \hat{\vartheta}_i) = \frac{1}{n} \sum_{j=1}^n x_j^k.$$

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

Použitelnost metody momentů

Možné problémy:

1. Potřebujeme existenci použitých momentů (zatímco jejich „odhady“ jsou definovány vždy).
2. \Rightarrow Nelze použít pro nenumerná data (např. volební předpověď), pokud je nelze smysluplně očíslovat.
3. Některé z rovnic nemusí být použitelné, protože nezávisí na parametrech nebo je nelze splnit. (Např. střední hodnota vychází teoreticky nulová, její odhad nulový být nemusí.)
Takové rovnice vyloučíme.

4. Snažíme se použít tolik rovnic (obvykle pro nejmenší možná k), abychom dostali konečný nenulový počet řešení.
5. Může být víc než jedno řešení (např. soustavy kvadratických rovnic).
6. Může být obtížné řešení nalézt.
7. Soustava může být špatně podmíněná (typicky pro velký počet parametrů).
8. Můžeme dospět k řešení, které **nesplňuje předpoklady**, $\hat{\boldsymbol{\theta}} \notin \Pi$ (např. parametry nemohou být libovolná čísla) \Rightarrow **vždy kontrolujte řešení!**
9. Všem rovnicím je přikládána stejná důležitost, což bývá nežádoucí (typicky pro velký počet parametrů), přitom i jejich fyzikální rozměry bývají různé.

Výhody:

1. Shoda momentů zajišťuje „podobné“ rozdělení modelu i dat.
2. Lze použít pro diskrétní, spojitě i **smíšené** rozdělení.

8.12.2 Metoda maximální věrohodnosti (angl. likelihood)

Motivace: Occamova břitva

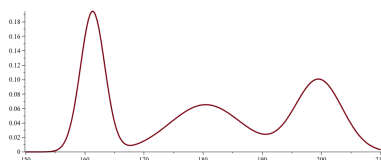
„Přikloňme se k nejjednoduššímu a nejpřirozenějšímu vysvětlení“

Příklad. 1000 respondentů vyjádřilo preference 3 politickým stranám následovně:

A	B	C
200	300	500

Strana A by měla uznat, že její preference jsou malé, místo aby výsledek přičítala statistické chybě nebo zfalšování průzkumu.

Příklad. V populaci nám vyšlo následující rozdělení parametru (např. tělesná výška):



Rozumný závěr je, že populaci tvoří 3 skupiny, v nichž je rozdělení podobné normálnímu (s různými středními hodnotami). Není vhodné usuzovat, že celé rozdělení je normální, nebo že ho tvoří 42 takových skupin.

Problém: Pro globální oteplování jsou podobné argumenty, ale ne tak jednoznačné.

Pravidlo: Z navržených modelů vybíráme ten, pro který jsou naměřené hodnoty „nejméně překvapivé.“

Pro diskrétní rozdělení

Pravděpodobnost realizace je funkce $L: \Pi \rightarrow \langle 0, 1 \rangle$, $\Pi \subseteq \mathbb{R}^i$, parametrů $\boldsymbol{\theta} = (\theta_1, \dots, \theta_i)$, zvaná **věrohodnost realizace diskrétního rozdělení**,

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

Hledáme takové hodnoty $\hat{\boldsymbol{\vartheta}} = (\hat{\vartheta}_1, \dots, \hat{\vartheta}_i)$, které maximalizují věrohodnost, resp. její logaritmus (*angl. 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.)

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.

Poznámka: Často určujeme věrohodnost až na násobek konstantou, kterou by bylo obtížné určit. Na výsledek optimalizace to nemá vliv.

Pro spojitě rozdělení

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

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

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. Může být více než jedno řešení.
(Může se stát, že různé hodnoty parametrů popisují totéž rozdělení – vadí to?)
2. Řešení nemusí existovat (pokud věrohodnostní funkce je nespojitá nebo parametrický prostor neuzavřený).
3. Může být obtížné řešení nalézt. (Uvznutí v lokálním extrému; parametrický prostor nemusí být souvislý.)
4. Hodnoty věrohodnosti mohou být velmi malé.
5. Fyzikální rozměr věrohodnosti může být ve spojitěm případě kuriózní.
6. **Nelze použít pro smíšené rozdělení!**

Výhody:

1. Hledání optima je o něco snazší než řešení soustavy rovnic – vždy nějakou aproximaci najdeme.
2. Různým datům je dán společný (srovnatelný) význam.
3. Lze použít i na nenumernická data.
4. Někdy vybíráme jen z konečného počtu modelů, pak metoda maximální věrohodnosti je použitelná, zatímco postupy založené na řešení rovnic nikoli.

8.12.3 Příklady na odhady parametrů

Cvičení. Odhadněte diskrétní rozdělení z četností hodnot v realizaci:

výsledek s	1	2	3	Σ
pravděpodobnost	a	b	c	1
četnost n_s	10	12	10	32

Řešení. *Metoda momentů:*

$$\begin{aligned}
 a + b + c &= 1, \\
 EX &= \sum_s s p_s = a + 2b + 3c = \frac{1}{n} \sum_s s n_s = \frac{10 + 2 \cdot 12 + 3 \cdot 10}{32} = 2, \\
 EX^2 &= \sum_s s^2 p_s = a + 4b + 9c = \frac{1}{n} \sum_s s^2 n_s = \frac{10 + 4 \cdot 12 + 9 \cdot 10}{32} = \frac{37}{8}, \\
 a = c &= \frac{5}{16}, \quad b = \frac{6}{16}.
 \end{aligned}$$

Metoda maximální věrohodnosti:

$$\begin{aligned}
 L(a, b) &= a^{10} \cdot b^{12} \cdot \underbrace{(1 - a - b)^{10}}_c, \\
 \ell(b) &= \ln \ell(a, b) = 10 \ln a + 12 \ln b + 10 \ln(1 - a - b), \\
 0 &= \frac{\partial}{\partial a} \ell(a, b) = \frac{10}{a} - \frac{10}{1 - a - b} = \frac{10}{a} - \frac{10}{c}, \\
 0 &= \frac{\partial}{\partial b} \ell(a, b) = \frac{12}{b} - \frac{10}{1 - a - b} = \frac{12}{b} - \frac{10}{c}, \\
 a &= c = \frac{5}{6} b,
 \end{aligned}$$

a z jednotkového součtu pravděpodobností opět

$$a = c = \frac{5}{16}, \quad b = \frac{6}{16}.$$

□

Obě metody vedly na empirické rozdělení. Tento výsledek není náhodný:

Věta. Pokud nějaká hodnota vektoru parametrů odpovídá empirickému rozdělení, pak je odhadem podle metody momentů i maximálně věrohodným odhadem.

Optimálních hodnot parametrů může být víc (a mohou vést na stejné rozdělení).

Důkaz. Označme u_1, \dots, u_i ($i \leq n$) všechny **různé** hodnoty, které se vyskytly v realizaci \mathbf{x} , n_s četnost a $r_s = n_s/n$ relativní četnost hodnoty u_s . Máme odhadnout pravděpodobnosti

$$q_s = p_X(u_s), \quad s = 1, \dots, i, \quad \sum_{s=1}^i q_s = 1.$$

Připomeňme, že empirické rozdělení $\text{Emp}(\mathbf{x})$ nabývá hodnot u_1, \dots, u_i s pravděpodobnostmi po řadě r_1, \dots, r_i .

Metoda maximální věrohodnosti:

$$\begin{aligned}
 L(q_1, \dots, q_i) &= \prod_{j=1}^n p_X(x_j) = \prod_{s=1}^i (p_X(u_s))^{n_s} = \prod_{s=1}^i q_s^{n_s}, \\
 \ell(q_1, \dots, q_i) &= \ln L(q_1, \dots, q_i) = \sum_{s=1}^i n_s \ln q_s.
 \end{aligned}$$

Použijeme metodu Lagrangeových multiplikátorů, tj. přičteme c -násobek podmínky jednotkového součtu neznámých a hledáme globální maximum funkce

$$h(c, q_1, \dots, q_i) = \sum_{s=1}^i n_s \ln q_s + c \left(1 - \sum_{s=1}^i q_s \right),$$

$$0 = \frac{\partial}{\partial q_s} h(c, q_1, \dots, q_i) = \frac{n_s}{q_s} - c.$$

Hodnota $\frac{n_s}{q_s} = c$ je nezávislá na $s \in \{1, \dots, i\}$. Určíme ji z podmínky $1 = \sum_{s=1}^i q_s = \frac{1}{c} \sum_{s=1}^i n_s = \frac{n}{c}$

$$\frac{n_s}{q_s} = c = n,$$

$$q_s = \frac{n_s}{n} = r_s$$

(empirické rozdělení).

Metoda momentů:

$$EX^k = \sum_{s=1}^i q_s u_s^k = \frac{1}{n} \sum_{j=1}^n x_j^k = \frac{1}{n} \sum_{s=1}^i n_s u_s^k = \sum_{s=1}^i r_s u_s^k.$$

Řešením je $q_s = r_s$ (empirické rozdělení). Je to jediné řešení, neboť matice soustavy

$$\begin{bmatrix} u_1^1 & u_2^1 & \cdots & u_i^1 \\ u_1^2 & u_2^2 & \cdots & u_i^2 \\ \vdots & \vdots & \ddots & \vdots \\ u_1^i & u_2^i & \cdots & u_i^i \end{bmatrix}$$

(tzv. Vandermondova matice) je regulární, právě když čísla u_1, \dots, u_i jsou navzájem různá. □

Cvičení. Odhadněte meze a, b spojitého rovnoměrného rozdělení z realizace

1. (3, 7, 5, 8, 1),
2. (0, 0, 0, 0, 5),

Řešení. **Metoda momentů:**

$$EX = \frac{a+b}{2} = \frac{1}{n} \sum_{s=1}^n x_s,$$

$$EX^2 = (EX)^2 + DX = \left(\frac{a+b}{2} \right)^2 + \underbrace{\frac{(b-a)^2}{12}}_{DX} = \frac{a^2 + ab + b^2}{3} = \frac{1}{n} \sum_{s=1}^n x_s^2.$$

1. *Soustava*

$$EX = \frac{a+b}{2} = \frac{24}{5},$$

$$EX^2 = \frac{a^2 + ab + b^2}{3} = \frac{148}{5},$$

má 2 řešení

$$a \doteq 0.36, \quad b \doteq 9.24,$$

$$a \doteq 9.24, \quad b \doteq 0.36$$

První řešení je jediné správné.

2. Soustava

$$\begin{aligned} EX &= \frac{a+b}{2} = 1, \\ EX^2 &= \frac{a^2 + ab + b^2}{3} = 5, \end{aligned}$$

má 2 řešení

$$\begin{aligned} a &\doteq -2.5, & b &\doteq 4.5, \\ a &\doteq 4.5, & b &\doteq -2.5. \end{aligned}$$

Žádné není správné, neboť $5 \notin \langle a, b \rangle$.

Metoda maximální věrohodnosti: Spojitá hustota je konstantní $\frac{1}{b-a}$ na intervalu $\langle a, b \rangle$.

$$L(a, b) = \prod_{j=1}^n \frac{1}{b-a} = \left(\frac{1}{b-a} \right)^n,$$

pokud $x_j \in \langle a, b \rangle$ pro všechna j ; jinak je nulová. Věrohodnost je maximální, pokud $b-a$ je minimální, tj.

$$a = \min_j x_j, \quad b = \max_j x_j.$$

1. $a = 1, \quad b = 8.$
2. $a = 0, \quad b = 5.$

Příklad. 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

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

Řešení:

$$\begin{aligned} \hat{\mu} &= \bar{\mathbf{x}}, \\ \hat{r} &= \frac{1}{n} \sum_{j=1}^n x_j^2 - \hat{\mu}^2 = \frac{1}{n} \sum_{j=1}^n x_j^2 - \bar{\mathbf{x}}^2 = \widehat{\sigma_X^2} = D \text{Emp}(\mathbf{x}). \end{aligned}$$

Metoda maximální věrohodnosti:

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

$$\begin{aligned}
\lambda(\mu, r) &= \ln \Lambda(\mu, r) = \frac{-1}{2r} \sum_{j=1}^n (x_j - \mu)^2 - \frac{n}{2} \ln r - \frac{n}{2} \ln 2\pi, \\
0 = \frac{\partial}{\partial \hat{\mu}} \lambda(\hat{\mu}, r) &= \frac{1}{r} \sum_{j=1}^n (x_j - \hat{\mu}) = \frac{1}{r} \left(\sum_{j=1}^n x_j - n \hat{\mu} \right) = \frac{n}{r} (\bar{x} - \hat{\mu}), \\
&\Rightarrow \hat{\mu} = \bar{x}, \\
0 = \frac{\partial}{\partial \hat{r}} \lambda(\hat{\mu}, \hat{r}) &= \frac{1}{2\hat{r}^2} \sum_{j=1}^n (x_j - \hat{\mu})^2 - \frac{n}{2\hat{r}} = \frac{n}{2\hat{r}^2} \left(\frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 - \hat{r} \right), \\
&\Rightarrow \hat{r} = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 = \widehat{\sigma_X^2} = \text{D Emp}(\mathbf{x}).
\end{aligned}$$

Motivační příklad (směs normálních rozdělení – pokračování):

Spojité rozdělení, jehož hustota má více maxim, aproximujeme směsí normálních,

$$\text{Mix}_{(c_1, \dots, c_m)}(\text{N}(\mu_1, \sigma_1^2), \dots, \text{N}(\mu_m, \sigma_m^2)),$$

s hustotou

$$\sum_{i=1}^m c_i \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(-\frac{(t - \mu_i)^2}{2\sigma_i^2}\right).$$

Vektor parametrů je $(c_1, \dots, c_m; \mu_1, \dots, \mu_m; \sigma_1, \dots, \sigma_m)$ s omezujícími podmínkami

$$\begin{aligned}
\sigma_i &\geq 0, & i &= 1, \dots, m, \\
0 \leq c_i &\leq 1, & i &= 1, \dots, m, \\
\sum_{i=1}^m c_i &= 1.
\end{aligned}$$

Pokus o řešení:

$$\begin{aligned}
f_X(t) &= \sum_{i=1}^m c_i f_{\text{N}(\mu_i, \sigma^2)}(t) = \sum_{i=1}^m c_i \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(t - \mu_i)^2}{2\sigma^2}\right), \\
\Lambda(\mu, \mathbf{c}) &= \prod_{j=1}^n f_X(x_j) = \prod_{j=1}^n \sum_{i=1}^m c_i f_{\text{N}(\mu_i, \sigma^2)}(x_j) = \\
&= \prod_{j=1}^n \left(\frac{1}{\sqrt{2\pi} \sigma} \sum_{i=1}^m c_i \exp\left(-\frac{(x_j - \mu_i)^2}{2\sigma^2}\right) \right) = \\
&= \left(\frac{1}{\sqrt{2\pi} \sigma} \right)^n \prod_{j=1}^n \sum_{i=1}^m c_i \exp\left(-\frac{(x_j - \mu_i)^2}{2\sigma^2}\right), \\
\lambda(\mu, \mathbf{c}) &= \sum_{j=1}^n \ln \sum_{i=1}^m c_i f_{\text{N}(\mu_i, \sigma^2)}(x_j) = \\
&= -n \ln(\sqrt{2\pi} \sigma) + \sum_{j=1}^n \ln \sum_{i=1}^m c_i \exp\left(-\frac{(x_j - \mu_i)^2}{2\sigma^2}\right).
\end{aligned}$$

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 i -té složce směsi popíšeme koeficientem $\alpha_{j,i} \in \langle 0, 1 \rangle$, přičemž

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

1. Zvolíme náhodně různé střední hodnoty složek směsi μ_i a nenulové koeficienty c_i , $i = 1, \dots, m$, splňující $\sum_{i=1}^m c_i = 1$.
- E. Stanovíme stupně příslušnosti

$$\alpha_{j,i} := \frac{c_i f_{N(\mu_i, \sigma^2)}(x_j)}{\sum_{i'=1}^m c_{i'} f_{N(\mu_{i'}, \sigma^2)}(x_j)} = \frac{c_i \exp\left(\frac{-(x_j - \mu_i)^2}{2\sigma^2}\right)}{\sum_{i'=1}^m \left(c_{i'} \exp\left(\frac{-(x_j - \mu_{i'})^2}{2\sigma^2}\right)\right)}$$

(jmenovatel je normalizační faktor).

M. Aktualizujeme koeficienty složek směsi

$$c_i := \frac{\sum_{j=1}^n \alpha_{j,i}}{\sum_{i'=1}^m \sum_{j=1}^n \alpha_{j,i'}} = \frac{1}{n} \sum_{j=1}^n \alpha_{j,i}$$

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

$$\mu_i := \frac{\sum_{j=1}^n \alpha_{j,i} x_j}{\sum_{j=1}^n \alpha_{j,i}} = \frac{\sum_{j=1}^n \alpha_{j,i} x_j}{n c_i}.$$

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 ukáзка EM algoritmu; lze jej snadno rozšířit na více dimenzí a jiné typy směsí. Použití pro parametry směrů rozdělení je typické, ne však jediné možné.

Problém: Uvznutí v lokálním extrému.

EM algoritmus rozšiřuje možnosti použití metody maximální věrohodnosti.

Použití empirického rozdělení $\text{Emp}(\mathbf{x})$ v odhadech

veličina	realizace odhadu	nestranný
EX	$E \text{Emp}(\mathbf{x}) = \frac{1}{n} \sum_i x_i = \bar{\mathbf{x}}$	+
EX^k	$E(\text{Emp}(\mathbf{x})^k) = \frac{1}{n} \sum_i x_i^k = m_k$	+
DX	$D \text{Emp}(\mathbf{x}) = \frac{1}{n} \sum_i (x_i - \bar{\mathbf{x}})^2$	–
	$\frac{n}{n-1} D \text{Emp}(\mathbf{x}) = \frac{1}{n-1} \sum_i (x_i - \bar{\mathbf{x}})^2 = s_{\mathbf{x}}^2$	+
σ_X	$\sigma_{\text{Emp}(\mathbf{x})} = \sqrt{\frac{1}{n} \sum_i (x_i - \bar{\mathbf{x}})^2}$	–
	$\sqrt{\frac{n}{n-1}} \sigma_{\text{Emp}(\mathbf{x})} = \sqrt{\frac{1}{n-1} \sum_i (x_i - \bar{\mathbf{x}})^2} = s_{\mathbf{x}}$	–
$q_X(\frac{1}{2})$	$q_{\text{Emp}(\mathbf{x})}(\frac{1}{2})$?
$q_X(\alpha)$	$q_{\text{Emp}(\mathbf{x})}(\alpha)$?
p_X bez omezení	$p_{\text{Emp}(\mathbf{x})}$	+
p_X s omezením	?	?
f_X	–	?

? = záleží na okolnostech

9 Testování hypotéz

9.1 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.

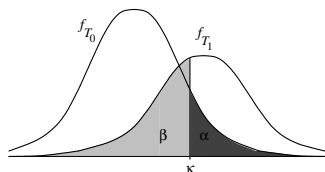
Příklad: Máme zastavit používání léku pro podezření z nežádoucích účinků?

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

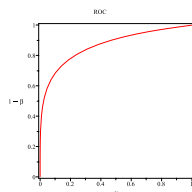
Alternativní hypotéza H_1 : Výrobce je viný, riziko se zvyšuje.

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).



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 rovnosti).

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 **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 .

Slovníček pojmů (pro porozumění jiným textům, zde se téměř nepoužijí)

test	skutečnost	anomální H_0 neplatí	normální H_0 platí	celkem
pozitivní, H_0 zamítnuta		TP	FP	P'
negativní, H_0 nezamítnuta		FN	TN	N'
celkem		P	N	

(Položky v tabulce mohou být pravděpodobnosti empirického nebo skutečného rozdělení nebo empirické četnosti.)

TP skutečně pozitivní (*true positive*)

FP falešně pozitivní, chyba 1. druhu (*false positive, type I error*)

TN skutečně negativní (*true negative*)

FN falešně negativní, chyba 2. druhu (*false negative, type II error*)

$$\alpha = \frac{FP}{N} = \frac{FP}{TN + FP} = \text{pravděpodobnost chyby 1. druhu}$$

$$\beta = \frac{FN}{P} = \frac{FN}{TP + FN} = \text{pravděpodobnost chyby 2. druhu}$$

$$\frac{TP}{TP + FN} = \frac{TP}{P} = 1 - \beta \quad \begin{array}{l} \text{senzitivita, síla, míra skutečně pozitivních} \\ \text{sensitivity, recall, true positive rate} \end{array}$$

$$\frac{TN}{TN + FP} = \frac{TN}{N} = 1 - \alpha \quad \begin{array}{l} \text{specificita, míra skutečně negativních} \\ \text{specificity, true negative rate} \end{array}$$

$$\frac{FP}{TN + FP} = \frac{FP}{N} = \alpha \quad \begin{array}{l} \text{míra falešně pozitivních} \\ \text{false positive rate} \end{array}$$

$$\frac{TP + TN}{TP + TN + FP + FN} = \frac{TP + TN}{P + N} \quad \begin{array}{l} \text{nesprávně přesnost} \\ \text{accuracy} \end{array}$$

$$\frac{TP + FN}{TP + TN + FP + FN} = \frac{P}{P + N} \quad \text{prevalence}$$

$$\frac{TP}{TP + FP} = \frac{TP}{P'} \quad \begin{array}{l} \text{přesnost,} \\ \text{prediktivní hodnota pozitivního testu} \\ \text{precision, positive predictive value} \end{array}$$

$$\frac{TN}{TN + FN} = \frac{TN}{N'} \quad \begin{array}{l} \text{prediktivní hodnota negativního testu} \\ \text{negative predictive value} \end{array}$$

Jednoduchá hypotéza: nulové hypotéze odpovídá jediná hodnota parametru.

Složená hypotéza: nulové hypotéze odpovídá více hodnot parametru.

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 nepadne 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: Pro mezní případ nulové hypotézy, $\vartheta = c$, odvodíme rozdělení testovací statistiky T , která s ϑ roste. Kvantily tohoto rozdělení určují intervalový odhad s koeficientem spolehlivosti $1 - \alpha$. Nulovou hypotézu zamítneme, pokud realizace t statistiky T padne mimo tento interval:

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:

H_0	H_1		H_0	H_1
$\vartheta = c$	$\vartheta > c$	nahradíme	$\vartheta \leq c$	$\vartheta > c$
$\vartheta = c$	$\vartheta < c$		$\vartheta \geq c$	$\vartheta < c$

9.2 Testy střední hodnoty normálního rozdělení

9.2.1 Při známém rozptylu σ^2

Výběrový průměr \bar{X}_n , který má rozdělení $N(\mu, \frac{\sigma^2}{n})$. S pravděpodobností $1 - \alpha$ je

$$q_{N(\mu, \sigma^2/n)}(\frac{\alpha}{2}) \leq \bar{X}_n \leq q_{N(\mu, \sigma^2/n)}(1 - \frac{\alpha}{2}),$$

po znormování

$$\Phi^{-1}(\frac{\alpha}{2}) = -\Phi^{-1}(1 - \frac{\alpha}{2}) \leq \text{norm } \bar{X}_n = \underbrace{\frac{\bar{X}_n - \mu}{\sigma} \sqrt{n}}_T \leq \Phi^{-1}(1 - \frac{\alpha}{2}),$$

V nulové hypotéze daná hodnota c nahrazuje neznámou střední hodnotu μ .

Realizaci t testové statistiky T ,

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

porovnáváme s kvantily **normovaného normálního rozdělení**:

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

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má nenulový rozptyl a rozsah výběru je velký.

9.2.2 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
$\mu = c$	$ t > q_{t(n-1)}(1 - \frac{\alpha}{2})$	$2(1 - F_{t(n-1)}(t))$
$\mu \leq c$	$t > q_{t(n-1)}(1 - \alpha)$	$1 - F_{t(n-1)}(t)$
$\mu \geq c$	$t < -q_{t(n-1)}(1 - \alpha)$	$F_{t(n-1)}(t)$

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má nenulový rozptyl a rozsah výběru je velký (pak můžeme místo Studentova rozdělení použít normální).

9.3 Testy rozptylu normálního rozdělení

S pravděpodobností $1 - \alpha$ je

$$q_{\chi^2(n-1)}(\frac{\alpha}{2}) \leq \underbrace{\frac{(n-1)S_x^2}{\sigma^2}}_T \leq q_{\chi^2(n-1)}(1 - \frac{\alpha}{2}),$$

V nulové hypotéze daná hodnota c nahrazuje neznámý rozptyl σ^2 .

Realizaci t testové statistiky T ,

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

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

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

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má nenulový rozptyl a rozsah výběru je velký (pak můžeme místo χ^2 -rozdělení použít normální).

9.4 Porovnání dvou normálních rozdělení

Předpoklad: **Nezávislé** výběry

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

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

9.4.1 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 $x > 0$:

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

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

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

$$\begin{aligned} U &:= \frac{(m-1)S_X^2}{\sigma^2} \text{ má } \chi^2(m-1), \\ V &:= \frac{(n-1)S_Y^2}{\sigma^2} \text{ má } \chi^2(n-1), \\ \xi &:= m-1, \eta := n-1, \\ F &= \frac{\frac{U}{\xi}}{\frac{V}{\eta}} = \frac{\frac{(m-1)S_X^2}{(m-1)\sigma^2}}{\frac{(n-1)S_Y^2}{(n-1)\sigma^2}} = \frac{S_X^2}{S_Y^2} = T. \end{aligned}$$

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 $\frac{S_X^2}{S_Y^2}$ 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
$DX = DY$	$t > q_{F(m-1, n-1)}(1 - \frac{\alpha}{2})$	$2(1 - F_{F(m-1, n-1)}(t))$

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

$$t' = \frac{1}{t} = \frac{s_y^2}{s_x^2} \geq 1$$

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

H_0	zamítáme pro	dosažená významnost
$DX \leq DY$	nezamítáme	žádná
$DX \geq DY$	$t' > q_{F(n-1, m-1)}(1-\alpha)$	$1 - F_{F(n-1, m-1)}(t')$
$DX = DY$	$t' > q_{F(n-1, m-1)}(1 - \frac{\alpha}{2})$	$2(1 - F_{F(n-1, m-1)}(t'))$

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má nenulový rozptyl a rozsah výběru je velký.

9.4.2 Testy středních hodnot dvou normálních rozdělení se stejným **známým** rozptylem σ^2

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

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 9.2.1).

9.4.3 Testy středních hodnot dvou normálních rozdělení s různými **známými** rozptily σ_X^2, σ_Y^2

$$\begin{aligned}\bar{X}_m &\text{ má } N(EX, \frac{\sigma_X^2}{m}), \\ \bar{Y}_n &\text{ má } N(EY, \frac{\sigma_Y^2}{n}), \\ \bar{X}_m - \bar{Y}_n &\text{ má } N(EX - EY, (\frac{\sigma_X^2}{m} + \frac{\sigma_Y^2}{n})).\end{aligned}$$

Za předpokladu $EX = EY$:

$$T := \frac{\bar{X}_m - \bar{Y}_n}{\sqrt{\frac{\sigma_X^2}{m} + \frac{\sigma_Y^2}{n}}} \text{ má } N(0, 1).$$

Testujeme realizaci t na $N(0, 1)$ (viz kapitola 9.2.1).

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má nenulový rozptyl a rozsah výběru je velký.

9.4.4 Testy středních hodnot dvou normálních rozdělení se stejným **neznámým** rozptylem σ^2

Nejprve ověříme předpoklad $DX = DY = \sigma^2$ (viz kapitola 9.4.1). *(Ve skutečnosti nemůžeme předpoklad ověřit, jedině vyvrátit; pokusíme se o to, a pokud se to nepodaří, pokračujeme. Bez tohoto předpokladu by byl další postup složitější, viz např. [Mood a kol.].)*

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

Za předpokladu $EX = EY$:

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

Pro měření významnosti potřebujeme odhad rozptylu. Máme dva odhady S_X^2, S_Y^2 téže hodnoty σ^2 ; použijeme jejich vážený průměr takový, abychom znali i jeho rozdělení.

$$\begin{aligned}\frac{(m-1)S_X^2}{\sigma^2} &\text{ má } \chi^2(m-1), \\ \frac{(n-1)S_Y^2}{\sigma^2} &\text{ má } \chi^2(n-1), \\ \frac{(m-1)S_X^2 + (n-1)S_Y^2}{\sigma^2} &\text{ má } \chi^2(m+n-2)\end{aligned}$$

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 , vedoucí na (vychýlený) odhad směrodatné odchylky

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

Ten použijeme místo neznámé směrodatné odchylky σ a výsledné kritérium

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

Testujeme realizaci t na rozdělení $t(m+n-2)$ (viz kapitola 9.2.2).

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má nenulový rozptyl a rozsah výběru je velký (pak můžeme místo Studentova rozdělení použít normální).

9.4.5 Testy středních hodnot dvou normálních rozdělení - párový test

(inspirováno [SH10], volně upraveno)

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 (vyjádřenou náhodnými veličinami Z_j) a projevuje se synchronně v obou výběrech; proto výběry **nelze popsat jako stejně rozložené a navzájem nezávislé**.

Situaci můžeme popsat následujícím modelem:

$$\begin{aligned} X_j &= Z_j + U_j, \\ Y_j &= Z_j + V_j - c, \end{aligned}$$

kde náhodné veličiny $U_1, \dots, U_n, V_1, \dots, V_n$ jsou nezávislé, U_1, \dots, U_n mají rozdělení $N(0, \sigma_U^2)$, V_1, \dots, V_n mají rozdělení $N(0, \sigma_V^2)$ a $c \in \mathbb{R}$. Testujeme hypotézu o hodnotě c (nejčastěji testujeme předpoklad $c = 0$).

Náhodné veličiny $\Delta_j := X_j - Y_j = U_j - V_j + c$ ($j = 1, \dots, n$) s rozdělením $N(c, \sigma_\Delta^2)$ (kde $\sigma_\Delta^2 := \sigma_U^2 + \sigma_V^2$) jsou nezávislé. Obecněji nám stačí:

Předpoklad: Náhodné veličiny $\Delta_j := X_j - Y_j$ jsou nezávislé a mají rozdělení $N(c, \sigma_\Delta^2)$.

Pokud rozptyl σ_Δ^2 známe, testujeme

$$T := \frac{\bar{\Delta} - c}{\sigma_\Delta} \sqrt{n} = \frac{\bar{\mathbf{X}} - \bar{\mathbf{Y}} - c}{\sigma_\Delta} \sqrt{n}$$

na rozdělení $N(0, 1)$ dle kapitoly 9.2.1.

Pokud rozptyl neznáme, testujeme

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

na rozdělení $t(n-1)$ dle kapitoly 9.2.2.

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má nenulový rozptyl a rozsah výběru je velký (pak můžeme místo Studentova rozdělení použít normální).

9.5 Korelace, její odhad a testování

(dle [Likeš, Machek])

Na základě realizace dvojrozměrného náhodného výběru $((x_1, y_1), \dots, (x_n, y_n))$ můžeme korelaci

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

odhadnout pomocí korelace empirického rozdělení neboli **realizace výběrového koeficientu korelace**

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

což je kosinus úhlu vektorů

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

Jednoduchý vzorec:

$$\begin{aligned} r_{\mathbf{x}, \mathbf{y}} &= \frac{n \sum_{j=1}^n x_j y_j - \left(\sum_{j=1}^n x_j\right) \left(\sum_{j=1}^n y_j\right)}{\sqrt{\left(n \sum_{j=1}^n x_j^2 - \left(\sum_{j=1}^n x_j\right)^2\right) \left(n \sum_{j=1}^n y_j^2 - \left(\sum_{j=1}^n y_j\right)^2\right)}} = \\ &= \frac{n}{n-1} \frac{\frac{1}{n} \sum_{j=1}^n x_j y_j - \bar{\mathbf{x}} \bar{\mathbf{y}}}{s_{\mathbf{x}} s_{\mathbf{y}}}. \end{aligned} \quad (2)$$

Výběrový koeficient korelace je odpovídající odhad

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

9.5.1 Test nekorelovanosti dvou **normálních** rozdělení

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

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 9.2.2 (pro oboustranný test i jednostranné testy).

Díky centrální limitní větě je odhad použitelný i pro výběr z jiného než normálního rozdělení, pokud má nenulový rozptyl a rozsah výběru je velký (pak můžeme místo Studentova rozdělení použít normální).

9.6 χ^2 -test dobré shody

9.6.1 Základní podoba testu

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

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

H_0 : Náhodná veličina má diskretní 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. **realizace četností** n_i (**empirické četnosti**) nebo **realizace relativních četností** $\frac{n_i}{n}$ ($i = 1, \dots, k$). Porovnáváme je s **teoretickými četnostmi** $n p_i$.

Speciální případ: Pro $k = 2$ mají N_1, N_2 binomická rozdělení $\text{Bi}(n, p_1)$, $\text{Bi}(n, p_2)$, která lze pro velká n přibližně nahradit normálními,

$$\begin{aligned} N(n p_1, n p_1 (1 - p_1)) &= N(n p_1, n p_1 p_2), \\ N(n p_2, n p_2 (1 - p_2)) &= N(n p_2, n p_1 p_2). \end{aligned}$$

Kvadráty normovaných veličin

$$(\text{norm } N_1)^2 = \frac{(N_1 - n p_1)^2}{n p_1 p_2}, \quad (\text{norm } N_2)^2 = \frac{(N_2 - n p_2)^2}{n p_1 p_2}$$

mají přibližně rozdělení $\chi^2(1)$. Jsou to tytéž náhodné veličiny, neboť

$$N_2 - n p_2 = n - N_1 - n (1 - p_1) = -(N_1 - n p_1),$$

takže rozdělení přibližně $\chi^2(1)$ má náhodná veličina

$$\begin{aligned} (\text{norm } N_1)^2 &= p_2 (\text{norm } N_1)^2 + p_1 (\text{norm } N_1)^2 \\ &= p_2 (\text{norm } N_1)^2 + p_1 (\text{norm } N_2)^2 \\ &= \frac{(N_1 - n p_1)^2}{n p_1} + \frac{(N_2 - n p_2)^2}{n p_2} = \sum_{i=1}^2 \frac{(N_i - n p_i)^2}{n p_i}. \end{aligned}$$

Obecně pro libovolné k je testovací statistikou

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

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

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

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

Modifikace: Pokud chceme naopak odhalit překvapivě dobrou shodu s modelem, použijeme dolní intervalový odhad a nulovou hypotézu zamítáme pro $t < q_{\chi^2(k-1)}(\alpha)$. Dosažená významnost: $F_{\chi^2(k-1)}(t)$.

Cvičení. 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 voličů, ale poslední má více než 2× méně poslanců).

9.6.2 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 (aspoň 5), aby náš předpoklad byl oprávněný.

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

Poznámka: Pokud data předpokládané rozdělení nemají a rozsah výběru zvětšíme $m \times$, systematický příspěvek ke kritériu se rovněž zvýší $\frac{m^2}{m} = m \times$. Proto se nedivme velkým hodnotám kritéria (velké síle testu) pro rozsáhlé výběry.

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.

9.6.3 χ^2 -test nezávislosti dvou rozdělení

(dle [Likeš, Machek])

H_0 : Dvě diskrétní 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 ; potřebujeme pouze četnosti N_{ij} , resp. jejich realizace 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 $p_i q_j$ ($i = 1, \dots, k$; $j = 1, \dots, m$),

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

Použijeme realizaci odhadu

$$t := \sum_{i=1}^k \sum_{j=1}^m \frac{(n_{ij} - n p_i q_j)^2}{n p_i q_j},$$

kde neznámé parametry p_i, q_j odhadneme pomocí maxima věrohodnosti neboli parametry empirického rozdělení,

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

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 $t > q_{\chi^2((k-1)(m-1))}(1-\alpha)$. Dosažená významnost: $1 - F_{\chi^2((k-1)(m-1))}(t)$.

9.6.4 χ^2 -test dobré shody dvou rozdělení

H_0 : Dva náhodné výběry pocházejí ze stejného diskrétního rozdělení.

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

Sjednocení obou výběrů považujeme za $m+n$ realizací náhodné veličiny X s neznámými pravděpodobnostmi p_i ($i = 1, \dots, k$), jejichž maximálně věrohodný odhad je

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

Zavedeme druhou náhodnou veličinu Y se dvěma hodnotami (např. 1, 2), které označují příslušnost k prvnímu, resp. druhému z původních výběrů. Např. z výběrů

$$(4, 5, 6, 4), \\ (6, 6, 6, 5, 5, 4)$$

vytvoříme dvojrozměrný výběr

$$((4, 1), (5, 1), (6, 1), (4, 1), (6, 2), (6, 2), (6, 2), (5, 2), (5, 2), (4, 2)).$$

Za předpokladu platnosti nulové hypotézy je X nezávislé na Y , což můžeme testovat stejně jako v předchozí metodě (na rozdělení $\chi^2(k-1)$).

Nulovou hypotézu zamítáme pro $t > q_{\chi^2(k-1)}(1-\alpha)$.

Dosažená významnost: $1 - F_{\chi^2(k-1)}(t)$.

Praktičtější (ekvivalentní) vzorec [Mood a kol.]:

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

9.7 Neparametrické testy

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

9.7.1 Znaménkový test

Rozlišujeme pouze znaménko odchylky od zvolené hodnoty c . Tím ztrácíme kvantitativní 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

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

(Podobně pro jednostranné testy.) Výpočet kvantilů je pracný, ale kritické hodnoty jsou tabelovány (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 testu.

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 $n \ln n$.

9.7.2 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.

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.
- [Kalina, Bacigál, Schiesslová] Kalina, M., Bacigál, T., Schiesslová, A.: *Základy pravděpodobnosti a matematické statistiky*. STU Bratislava, 2010.
- [Kalina, Minarechová] Kalina, M., Minarechová, Z.: *Applied Mathematics For Civil Engineers*. STU Bratislava, 2015.
- [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.
- [Gonick & Smith] Gonick, L., Smith, W.: *The Cartoon Guide to Statistics*. HarperPerennial 1993, HarperCollins Publishers, New York, 2005.
- [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.