

## Přednáška 08a

### Od popisné statistiky k matematické statistice

Statistický soubor ... množina statistických jednotek

Měříme jeden nebo více znaků

Typy znaků - nominální ( $x_1 = x_2?$ )

- ordinální ( $x_1 < x_2?$ )

- intervalové ... jsme navíc schopni prohodit  
řadí  $x_1 - x_2$

- poměrové ... nerovnost, řadí, řadí

Nechtě  $x_1, x_2, \dots, x_n$  je soubor hodnot, které lze uspořádat, namísto měření na  $n$  statistických jednotkách.

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

Různé hodnoty nechtě jsou  $a_1 < a_2 < \dots < a_m$ ,  
 $n_i$  je počet jednotek s naměřenou hodnotou  $a_i$   
(lůžní četnost).

### Průměry

① Aritmetický průměr

$$\bar{x} = \frac{1}{n} \left( \sum_{i=1}^n x_i \right) = \frac{1}{n} \left( \sum_{j=1}^m n_j a_j \right)$$

Aritmetický průměr je invariantní vůči afinnímu transformaci

$$y_i = ax_i + b$$

Pak

$$\bar{y} = a\bar{x} + b$$

② Geometrický průměr pro kladná čísla

$$\bar{x}^G = \sqrt[n]{x_1 x_2 \dots x_n}$$

③ Harmonický průměr pro kladná čísla

$$\bar{x}^H = \left\{ \frac{1}{n} \left( \sum_{i=1}^n \frac{1}{x_i} \right) \right\}^{-1}$$

Platí

$$\bar{x}^H \leq \bar{x}^G \leq \bar{x}$$

### Median, kvartil, percentil

Median souboru  $x_1, x_2, \dots, x_n$  je číslo  $\tilde{x}$  takové, že 50% hodnot  $x_1, x_2, \dots, x_n$  je  $\leq \tilde{x}$  a zbytek hodnoty jsou  $> \tilde{x}$ . Tímto není  $\tilde{x}$  nutně jednoznačné.

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

n liché

$$\tilde{x} = x_{\left(\frac{n+1}{2}\right)}$$

n sudé

$$\tilde{x} = \frac{x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n}{2}+1\right)}}{2}$$

Pro  $0 < p < 1$  je  $p$ -tý kvantil,  $p$ -tý percentil číslo  $x_p$  takové, že  $p\%$  hodnot je  $\leq x_p$  a zbylé jsou  $> x_p$ . Nemá-li takové číslo, bereme ho jako aritmetický průměr dvou po sobě jdoucích

Hodnoty  $x_i$       1 1 2 2 3 3 3 3

$$x_{0,5} = \tilde{x} = \frac{2+3}{2} = 2,5$$

$$x_{0,25} = Q_1 = \text{dolní kvantil je } \frac{1+2}{2} = 1,5$$

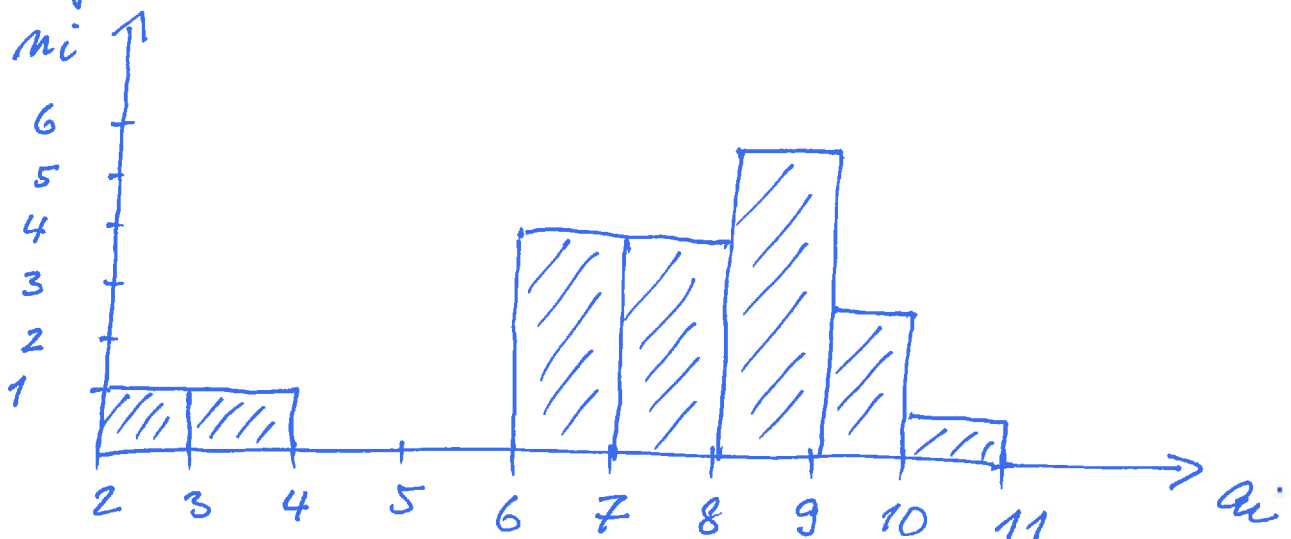
$$x_{0,75} = Q_3 = \text{horní kvantil je } \frac{3+3}{2} = 3$$

Modus je hodnota s největší četností  $\hat{x} = 3$ .

### Příklad

$a_i$	3	4	7	8	9	10	11
$m_i$	1	1	4	4	6	3	1

### Histogram

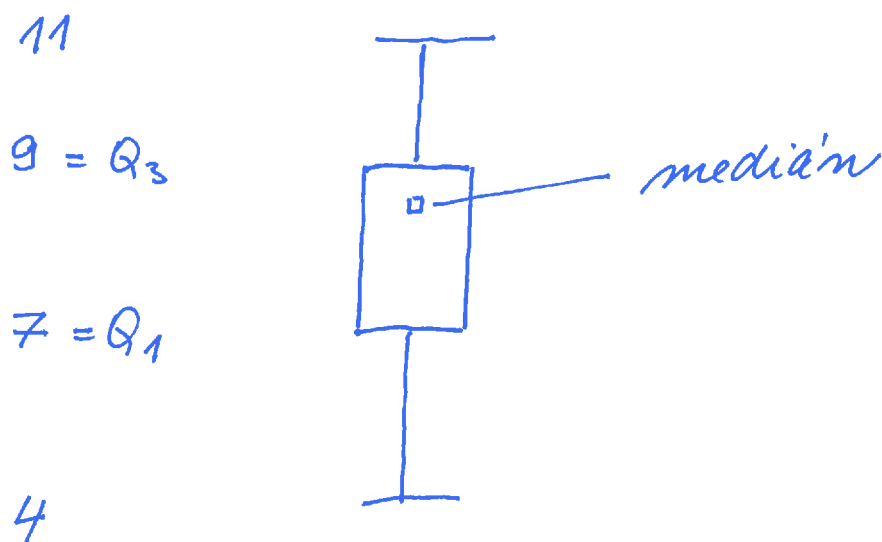


### Kračový graf

$$n = 20, \quad \bar{x} = 8,1, \quad \tilde{x} = \frac{x_{(10)} + x_{(11)}}{2} = \frac{8+9}{2} = 8,5$$

$$x_{0,25} = Q_1 = \frac{x_{(5)} + x_{(6)}}{2} = 7$$

$$x_{0,75} = Q_3 = \frac{x_{(15)} + x_{(16)}}{2} = 9$$



### Rozptyl a směrodatná odchylka

Rozptyl 
$$s_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$s_x \geq 0$  je směrodatná odchylka

### Průměrná odchylka

$$D_x = \frac{1}{n} \sum_{i=1}^n |x_i - \tilde{x}|$$

Věta (a) Funkce  $S(t) = \frac{1}{n} \sum_{i=1}^n (x_i - t)^2$  nabývá svého minima pro  $t = \bar{x}$ .

(b) Funkce  $D(t) = \frac{1}{n} \sum_{i=1}^n |x_i - t|$  nabývá svého minima pro  $t = \tilde{x}$ .

Důkaz (a)

$$\begin{aligned} \sum_{i=1}^n (x_i - t)^2 &= \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - t)^2 = \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n (\bar{x} - t)^2 + 2 \sum_{i=1}^n (x_i - \bar{x})(\bar{x} - t) \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n (\bar{x} - t)^2 + 2 \left( \underbrace{\sum_{i=1}^n x_i - n\bar{x}}_0 \right) (\bar{x} - t) \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 + n (\bar{x} - t)^2 \end{aligned}$$

Odkud je vidět, že výraz má vždy své minimum pro  $t = \bar{x}$ .

(b) Zkoumáme podobně rozích

$$|x_{(1)} - t| + |x_{(n)} - t|$$

$$|x_{(2)} - t| + |x_{(n-1)} - t|$$

atd

Pokud

$$|x_{(1)} - t| + |x_{(n)} - t| = |x_{(1)} - \tilde{x}| + |x_{(n)} - \tilde{x}| \stackrel{\text{pro } t \in [x_{(1)}, x_{(n)}]}{=} x_{(n)} - x_{(1)}$$

$$\geq |x_{(1)} - \tilde{x}| + |x_{(n)} - \tilde{x}| \stackrel{\text{pro } t \notin [x_{(1)}, x_{(n)}]}{>} x_{(n)} - x_{(1)}$$

$$|x_{(2)} - t| + |x_{(n-1)} - t| = |x_{(2)} - \tilde{x}| + |x_{(n-1)} - \tilde{x}| = x_{(n-1)} - x_{(2)} \stackrel{\text{pro } t \in [x_{(2)}, x_{(n-1)}]}{=}$$

$$a \quad |x_{(2)} - t| + |x_{(n-1)} - t| > |x_{(2)} - \tilde{x}| + |x_{(n-1)} - \tilde{x}|$$

pre  $t \notin [x_{(2)}, x_{(n-1)}]$ .

ald. Odhad dokážeme, že  $D(t)$  nalyza' sne'ke minima pre  $t = \tilde{x}$ .

### Nomina'lni' analyza'

$n_i$  počet jednotiek s nomina'lnym znacenim  $a_i$

$a_1$	$a_2$	$a_3$	...	$a_k$
$n_1$	$n_2$	$n_3$		$n_k$

$$\sum_{i=1}^k n_i = n$$

Mera variability analyzy vyjadruje entropie

$$H_x = - \sum_{i=1}^k \frac{n_i}{n} \ln \left( \frac{n_i}{n} \right)$$

$$k=1 \quad H_x = 0$$

$$k=2 \quad n_1 = n_2$$

$$H_x = - \sum_{i=1}^2 \frac{1}{2} \ln \frac{1}{2} = \ln 2 > 0.$$

Entropie nalyza' sne'ke maxima pri stejne' cistoti analyzy:

Hledáme maximum funkce

$$H(p_1, \dots, p_k) = - \sum_{i=1}^k p_i \ln p_i$$

-7-

Pa podmínky  $\sum_{i=1}^k p_i = 1$  (a samozřejmě  $p_i > 0$ ). Pomocí Lagrangeova multiplikátoru  $\lambda$

$$\frac{\partial H}{\partial p_i} = -\ln p_i - 1 = \lambda$$

$$p_i = e^{\lambda-1}$$

$$k e^{\lambda-1} = 1$$

$$\lambda = 1 - \ln k$$

$$p_i = \frac{1}{k}$$

V bodě  $[\frac{1}{k}, \frac{1}{k}, \dots, \frac{1}{k}]$  má  $H$  svého lokálního maxima. Funkci  $f(y) = -y \ln y$  lze nejlehčeji rozšířit na  $y=0$  tak, aby  $\lim_{y \rightarrow 0^+} (-y \ln y) = 0$ .

$H$  má na kompaktní množině

$$\left\{ \sum p_i = 1, \quad p_i \geq 0 \right\}$$

svoje maxima, a to musí být v bodě  $[\frac{1}{k}, \frac{1}{k}, \dots, \frac{1}{k}]$ .

## Pravděpodobnost

Hlási me korekci - máme výsledky pro 1, 2, 3, 4, 5, 6.

Můžeme mít situace, kdy mají výsledků je nekonečně mnoho:

Operace body mincí - rub a kč u jednoho hodu.

$w_k \in \mathbb{N} \cup \{\infty\}$  kč padne poprvé v  $k$ -tém hodu

Pravděpodobnost

$$P(\omega_k) = \frac{1}{2^k}$$

$$P(\omega_\infty) = 0$$

$$\sum_{k=1}^{\infty} P(\omega_k) = 1.$$

### Jerové pole - $\sigma$ -algebra

Pracujeme se sákladní množinou  $\Omega$ . Vybrané podmnožiny reprezentují jevy.

Jerové pole na  $\Omega$  je systém podmnožin  $\mathcal{A}$  množiny  $\Omega$  splňující, že

- $\Omega \in \mathcal{A}$
- $A, B \in \mathcal{A} \Rightarrow A \setminus B \in \mathcal{A}$
- $A_i \in \mathcal{A}, I$  nejryšle spočetná  $\Rightarrow \bigcup_{i \in I} A_i \in \mathcal{A}$

$\mathcal{A}$  nazýváme také  $\sigma$ -algebrou.

### Důsledky:

- $\emptyset \in \mathcal{A}$  (nemohou být,  $\Omega$  je jistý jev)
- $A \in \mathcal{A}, A^c = \Omega \setminus A \in \mathcal{A}$  (doplnek)
- $A, B \in \mathcal{A} \Rightarrow A \cap B = A \setminus (\Omega \setminus B) \in \mathcal{A}$

### Pravděpodobnostní prostor

je jové pole  $\mathcal{A}$  společně s funkcí  $P: \mathcal{A} \rightarrow [0,1]$  splňující, že

$$(1) \quad P(\Omega) = 1,$$

(2)  $A_i$  jsou dvou disjunktní,  $I$  nejryšle spočetná



$$P\left(\bigcup_{i \in I} A_i\right) = \sum_{i \in I} P(A_i)$$

Důsledky:

•  $P(A^c) = 1 - P(A)$ ,  $P(\emptyset) = 0$ .

•  $P\left(\bigcup_{i=1}^k A_i\right) = \sum_{i=1}^k P(A_i) - \sum_{1 \leq i < j \leq k} P(A_i \cap A_j)$

+  $\sum_{1 \leq i < j < l \leq k} P(A_i \cap A_j \cap A_l) - \dots$

Věta

(1) Necht'  $A_1 \subseteq A_2 \subseteq \dots \subseteq A_n \subseteq \dots \in \mathcal{A}$

Pak  $P\left(\bigcup_{i=1}^{\infty} A_i\right) = \lim_{n \rightarrow \infty} P(A_n)$

(2) Necht'  $A_1 \supseteq A_2 \supseteq A_3 \supseteq \dots \in \mathcal{A}$

Pak  $P\left(\bigcap_{i=1}^{\infty} A_i\right) = \lim_{n \rightarrow \infty} P(A_n)$ .

Důkaz (1)  $A = \bigcup_{i=1}^{\infty} A_i = A_1 \cup (A_2 \setminus A_1) \cup (A_3 \setminus A_2) \cup \dots$

Podob  
 $P(A) = P(A_1) + \sum_{i=1}^{\infty} P(A_{i+1} \setminus A_i) = P(A_1) + \lim_{n \rightarrow \infty} \sum_{i=1}^n P(A_{i+1} \setminus A_i)$   
 $= P(A_1) + \lim_{n \rightarrow \infty} \{ P(A_2) - P(A_1) + P(A_3) - P(A_2) + \dots + P(A_{n+1}) - P(A_n) \}$

$$= P(A_1) + \lim_{n \rightarrow \infty} (P(A_{n+1}) - P(A_1)) = \lim_{n \rightarrow \infty} P(A_{n+1}).$$

$$\begin{aligned} (2) \quad P\left(\bigcap_{i=1}^{\infty} A_i\right) &= 1 - P\left(\bigcup_{i=1}^{\infty} (\Omega \setminus A_i)\right) = \\ &= 1 - \lim_{n \rightarrow \infty} P(\Omega \setminus A_n) = 1 - 1 + \lim_{n \rightarrow \infty} P(A_n) \\ &= \lim_{n \rightarrow \infty} P(A_n) \end{aligned}$$

V druhé rovnici jsme použili (1).

### Nesámkost jeví A a B

jevy A a B jsou nezávislé, tj. splňují

$$P(A \cap B) = P(A) \cdot P(B)$$

Základní předpoklad o pravidelném počtu:  
 $\Omega$  konečná množina, v níž je všech jevích podmnožin a

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

## Podmiňena' pravdĕpodobnost

$H$  je jím s nenulovou pravdĕpodobností.

Podmiňena' pravdĕpodobnost jím  $A$  sa podmiňky  $H$  je pravdĕpodobnost

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

Lemma Nechtĕ  $B$  je disj. sřednocením jím  $B_1 \cup B_2 \cup \dots \cup B_n$ . Pak

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

Speciálně ma  $\Omega = B_1 \cup B_2 \cup \dots \cup B_n$  disj.

$$P(A) = \sum P(A|B_i) P(B_i).$$

Důkaz:

$$P(A|B) = \frac{P(A \cap (B_1 \cup \dots \cup B_n))}{P(B_1 \cup \dots \cup B_n)} = \sum_{i=1}^n \frac{P(A \cap B_i)}{P(B)}$$

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

## Bayesova věta

(A) Bayesův vzorec pro inverzní pravděpodobnost

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

(B) 1. Bayesův vzorec

$$P(A|B) = \frac{P(A) p(B|A)}{p(A) p(B|A) + p(A^c) p(B|A^c)}$$

(C) obecnější verze: Je-li  $\Omega = A_1 \cup A_2 \cup \dots \cup A_n$  disjunktivně, pak

$$P(A_i|B) = \frac{p(A_i) p(B|A_i)}{\sum_{j=1}^n p(A_j) p(B|A_j)}$$

Důkaz: (A)

$$\begin{aligned} \frac{p(A) p(B|A)}{p(B)} &= \frac{p(A) \frac{p(B \cap A)}{p(A)}}{P(B)} = \frac{P(A \cap B)}{P(B)} \\ &= P(A|B) \end{aligned}$$

(B) a (C) za  $P(B)$  dosadíme z lemmatu.

Příklad v předmětu X je úspěšnost 40%,  
v předmětu Y je 80%. Každý předmět  
v daný den dělal nejvýš pět studentů.

žalá' je marděpodobnost, že student, který uspěl,  
dělal předmět X?

A student dělal předmět X

B student u slovníky uspěl

$A^c$  student dělal předmět Y

Víme, že  $P(B|A) = 0,4$ ,  $P(B|A^c) = 0,8$ .

Dále  $P(A) = P(A^c) = 0,5$ . Chceme spočítat  
 $P(A|B)$ . Dosazením do vzorce

$$\begin{aligned} P(A|B) &= \frac{p(A) \cdot p(B|A)}{p(A) p(B|A) + p(A^c) p(B|A^c)} \\ &= \frac{0,5 \cdot 0,4}{0,5 \cdot 0,4 + 0,5 \cdot 0,8} = \frac{1}{3} \end{aligned}$$



Nesánilost jeví<sup>o</sup> A a B pro  $P(B) \neq 0$  je ekvivalentní  
s rovností

$$P(A) = P(A|B).$$