

# STATISTIKA V KOSTCE

## DRUHÝ TEST

### Pravděpodobnost

Nespojité náhodná veličina:

Pravděpodobnostní funkce ( $P(x)$ ) a Distribuční funkce ( $F(x)$ ).

Spojité náhodná veličina: **(integrály)**

Distribuční funkce ( $F(x)$ ) a hustota pravděpodobnosti ( $f(x)$ ).

dolní hranice      horní hranice

$$P(x_1 < X \leq x_2) = P\left(\frac{x_1 - \mu}{\sigma} < \frac{X - \mu}{\sigma} \leq \frac{x_2 - \mu}{\sigma}\right) = P(u_1 < U \leq u_2) = \underline{\Phi(u_2)} - \underline{\Phi(u_1)}$$

Binomické rozdělení  $Bi(n, \pi)$

$$P(x) = \binom{n}{x} \pi^x (1-\pi)^{n-x}$$

*s vrácením, hazardní hrací kostka, karty, ruleta*

Hypergeometrické rozdělení  $Hg(M, N, n)$

$$P(x) = \frac{\binom{M}{x} \binom{N-M}{n-x}}{\binom{N}{n}}$$

*Bez vrácení - Kontrola, sarka.*

Vzadání předpoklad normálního rozdělení

*! hledáme v tabulkách*

Tabulka III.  
Distribuční funkce normálního rozdělení

### Matematická statistika

a)  $\sigma^2$  známý

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

*všestranná odchylka*

*smerodatná odchylka*

Vzadání  $\Rightarrow$  Určete interval spolehlivosti (pr. 90%; 95%; 99%; ...)

*zhádno odchylka (je zadána; obecně dana)*

b)  $\sigma^2$  neznámý

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

*Studentovo rozdělení*

*nezádano odchylka (masí se spočítat; vzadání je že byla spočítaná.)*

*pocet stupnic volnosti*

→ Tabulka VI. Kvantily rozdělení  $t$   $T \sim t(n-1)$

$$P\left(P - u_{1-\alpha/2} \sqrt{\frac{P(1-P)}{n}} < \pi < P + u_{1-\alpha/2} \sqrt{\frac{P(1-P)}{n}}\right) = 1 - \alpha$$

*zadáno relativně*

*normální rozdělení*

→ Tabulka IV. Kvantily normovaného normálního rozdělení ( $u_P$ )

### Testování statistických hypotéz

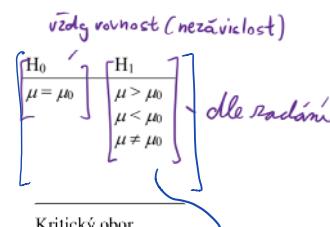
Jestliže hodnota testové statistiky:

Patří (nerovnost platí) do kritické oboru, potom na dané hladině významnosti **ZAMÍTÁME H0 a přijímáme H1**.

Nepatří (nerovnost neplatí) do kritického oboru, potom na dané hladině významnosti **NEZAMÍTÁME H0. Výsledek je statisticky nevýznamný**. Pozor: Neznamená to, že zamítáme H1.

+ řídí určidíme hladinu významnosti

(na 5% hl.významnosti; na 1% hl.významnosti)



Kritický obor

$$W_\alpha = \{u; u \geq u_{1-\alpha}\}$$

$$W_\alpha = \{u; u \leq -u_{1-\alpha}\}$$

$$W_\alpha = \{u; |u| \geq u_{1-\alpha/2}\}$$

olej H1 vybereme

kritický obor.

*známá/neznámá - stejný princip jako u intervalů*

$\sigma^2$  známý

$$U = \frac{\bar{X} - \mu_0}{\sigma} \sqrt{n} \quad U \sim N(0,1)$$

$$\sigma^2 \text{ neznámý} \quad \text{neznámé rozptyl a } n \leq 30$$

$$T = \frac{\bar{X} - \mu_0}{S} \sqrt{n} \quad T \sim t(n-1)$$

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

Testové kritérium

$$\sigma^2 \text{ neznámý } (n > 30)$$

$$U = \frac{\bar{X} - \mu_0}{S} \sqrt{n} \quad U \approx N(0,1)$$

Testové kritérium *zadáno relativně (%)*

$$U = \frac{P - \pi_0}{\sqrt{\pi_0(1-\pi_0)}} \quad U \approx N(0,1)$$



STATISTICKY  
NEKLASICKY

DOUČOVÁNÍ STATISTIKY S ADRIANOУ

# STATISTIKA V KOSTCE

## DRUHÝ TEST

Rovnost středních hodnot dvou rozdělení  
velké nezávislé výběry

(2 skupiny nebo 2 lektorání) , vždy v Excelu.

Testové kritérium

$$\sigma_1^2 \text{ a } \sigma_2^2 \text{ neznámé} \\ U = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \quad U \approx N(0,1)$$

nezávislé výběry (ženy/máci, 1. skupina/2. skupina)

váčné nebo Excel

EXCEL: Data - Analýza dat -  
Dvouvýběrový T-test  
s nerovností rozptylu

párový t-test: (stejná skupina)  
2. lektorání

Testové kritérium

$$T = \frac{\sqrt{n}\bar{D}}{S_D} \quad T \sim t(n-1)$$

závislé výběry (Před/Po,  
studenti prošli 1. a 2. testem)

### Chi-kvadrát test dobré shody

$$G = \sum_{j=1}^k \frac{(n_j - n\pi_{0,j})^2}{n\pi_{0,j}} \quad G \approx \chi^2(k-1)$$

\ počet stupňů volnosti

• více % a více hodnot

Kvantily rozdělení  $\chi^2$  - vždy tab. 7

DOMÁCÍ ÚKOL: Marketingový plán tvrdil, že záznam koncertu skupiny Kapička se prodá v poměru 70% CD, 20% DVD a 10% kazety. Za měsíc se skutečně prodá 2400 kusů CD, 1075 DVD a 384 kusů kazet tohoto koncertu. Ověřte, zda byl předpoklad marketingového plánu správný.

$$\begin{bmatrix} 70\% ; 20\% ; 10\% \Rightarrow \text{předpoklad} \\ 2400 ; 1075 ; 384 \Rightarrow \text{konkrétní počet} \end{bmatrix}$$

$$n = 2400 + 1075 + 384$$

### Analýza závislostí

Kontingenční tabulka ( $r \times s$ )

$$G = \sum_{i=1}^r \sum_{j=1}^s \frac{(n_{ij} - n_{i+}n_{+j}/n)^2}{n_{i+}n_{+j}/n} \quad G \approx \chi^2((r-1)(s-1))$$

$$G = \frac{n(n_{11}n_{22} - n_{12}n_{21})^2}{n_{1+}n_{+1}n_{2+}n_{+2}} \quad - 2 \times 2 \text{ tabulka}$$

$$\chi^2 = \sqrt{\frac{G}{n(m-1)}} \quad m = \min(r, s) \quad - \text{Cramérův koeficient} \quad \text{koeficient } \in [0, 1]$$

⇒ více skupin a více hodnot:

	A	B	C	D	E	$\Sigma$
máč	12	5	7	1	3	28
žena	3	8	4	2	1	18
$\Sigma$	15	13	11	3	2	46

$$C = \sqrt{\frac{G}{G+n}} \quad - \text{Pearsonův koeficient} \quad (\text{vzětý nízky rozdíly})$$

### Analýza rozptylu

$$F = \frac{\frac{S_{y.m}}{k-1}}{\frac{S_{y.v}}{n-k}} \Rightarrow \text{Excel}$$

dítě	2	5	3	8
máč	12	5	2	1
žena	3	8	4	2

více skupin, ale už se dál nečlení; jen kvádrum skupin.

$$P^2 = \frac{S_{y.m}}{S_y}$$

(Poměr determinace)  
síla závislosti



STATISTICKY  
NEKLASICKY

DOUČOVÁNÍ STATISTIKY S ADRIANOU

# STATISTIKA V KOSTCE

## DRUHÝ TEST

### Regres a korelace

regresní přímka  $y = \beta_0 + \beta_1 x + \varepsilon$ ,

regresní parabola:  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2$

vícenásobná regrese:  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_k x_k$

Zabýváme se dvojicí veličin:

Y - vysvětlovaná, závisle proměnná - vzdg 18

X - vysvětlující, nezávisle proměnná - 1 nebo více

*přímka  
vícenásobná  
nebo parabola*

$$R^2 = I^2 = \frac{S_T}{S_y} - \text{koeficient determinace } \in \langle 0,1 \rangle$$

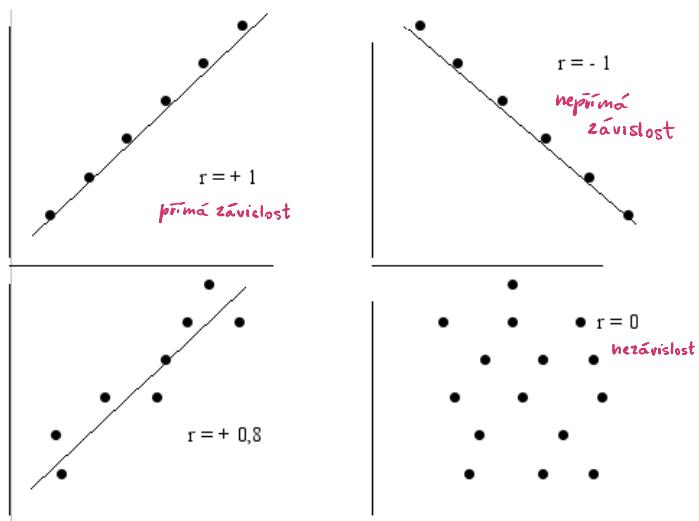
- požívá se pro porovnání více modelů.

$$I^2_{ADJ} = R^2_{ADJ} = 1 - (1 - I^2) \frac{n-1}{n-p} - \text{upravený koeficient determinace}$$

Je lepší přímka nebo parabola?

korelační koeficient

$$r_{yx} = r_{xy} = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{\sqrt{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \sqrt{n \sum_{i=1}^n y_i^2 - (\sum_{i=1}^n y_i)^2}} = \frac{\bar{xy} - \bar{x} \bar{y}}{\sqrt{(\bar{x}^2 - \bar{x}^2)(\bar{y}^2 - \bar{y}^2)}} = \frac{s_{xy}}{s_x s_y} = \sqrt{R^2} \quad r_{yx} \in \langle -1,1 \rangle$$



Zdroj: [absolventi.gymcheb.cz](http://absolventi.gymcheb.cz)

#### Regresní statistika

Násobné R	0,984743	- Korelační Koeficient
Hodnota spolehlivosti R	0,969718	- Koeficient determinace
Nastavená hodnota spolehlivosti	0,965933	- upravený koef. determinace
Chyba stř. hodnoty	58,59154	
Pozorování	10	

#### Test hypotézy o regresním parametru

	Koeficienty ba stř. hodn.	t Stat	Hodnota P	Dolní 95%	Horní 95%	Dolní 99,0%
Hranice	B <sub>0</sub> -160,347	41,00253	-3,91066	0,004477	-254,899	-65,7949
Soubor X 1	B <sub>1</sub> 7,573698	0,473188	16,00567	2,33E-07	6,482524	8,664873

Excel: Data - Analýza Dat - Regrese (x - nezávislá, y - závislá)

veřejně

#### Test o modelu

ANOVA

	Rozdíl	SS	MS	F	úznamnost F
Regres	1	879463,2	879463,2	256,1815	2,33E-07
Rezidua	8	27463,75	3432,969		
Celkem	9	906926,9			



STATISTICKY  
NEKLASICKY

DOUČOVÁNÍ STATISTIKY S ADRIANOU