

# Statistika

(MD360P03Z, MD360P03U)  
ak. rok 2007/2008

Karel Zvára

karel.zvara@mff.cuni.cz  
http://www.karlin.mff.cuni.cz/~zvára

6. listopadu 2007



Statistika (MD360P03Z, MD360P03U) ak. rok 2007/2008

## chování výběrového průměru

- ▶ necht  $X_1, X_2, \dots, X_n$  jsou nezávislé náhodné veličiny s **libovolným stejným rozdělením** se střední hodnotou  $\mu$  a rozptylem  $\sigma^2$ , tj. **náhodný výběr** z onoho rozdělení

- ▶ **průměr**  $X_1, X_2, \dots, X_n$ :

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

- ▶ připomeňme vlastnosti střední hodnoty ▶ Vlastnosti

$$\mu_{X+Y} = \mu_X + \mu_Y, \quad \mu_{b \cdot X} = b \cdot \mu_X$$

- ▶ proto je

$$\mu_{\bar{X}} = \mu_{\frac{1}{n} \cdot \sum_{i=1}^n X_i} = \frac{1}{n} \cdot \mu_{\sum_{i=1}^n X_i} = \frac{1}{n} \sum_{i=1}^n \mu_{X_i} = \frac{1}{n} n \mu = \mu$$

- ▶  $\mu_{\bar{X}} = \mu$ , tj.  $\bar{X}$  je **nestranný odhad** parametru  $\mu$

## populace a výběr

- ▶ populaci charakterizujeme pomocí parametrů rozdělení, případně typu rozdělení
- ▶ výsledek měření na náhodně vybraném prvku populace – náhodná veličina
- ▶ skutečné hodnoty parametrů neznáme
  - ▶ chceme je odhadnout
  - ▶ chceme rozhodnout o platnosti tvrzení (hypotézy) o parametrech
- ▶ jako výběr si představujeme několik **nezávislých** náhodných veličin se stejným rozdělením a neznámými parametry
  - ▶ parametry odhadujeme na základě výběru
  - ▶ o hypotézách rozhodujeme na základě výběru
- ▶ příklady
  - ▶ střední hodnotu náhodné veličiny (populační průměr) odhadujeme pomocí výběrového průměru
  - ▶ rozptyl náhodné veličiny odhadujeme pomocí výběrového rozptylu

6. přednáška 5. listopadu 2007

Statistika (MD360P03Z, MD360P03U) ak. rok 2007/2008

## variabilita výběrového průměru

- ▶ pro rozptyl **nezávislých** náhodných veličin platí ▶ Vlastnosti

$$\sigma_{X+Y}^2 = \sigma_X^2 + \sigma_Y^2 \quad \sigma_{b \cdot X}^2 = b^2 \sigma_X^2$$

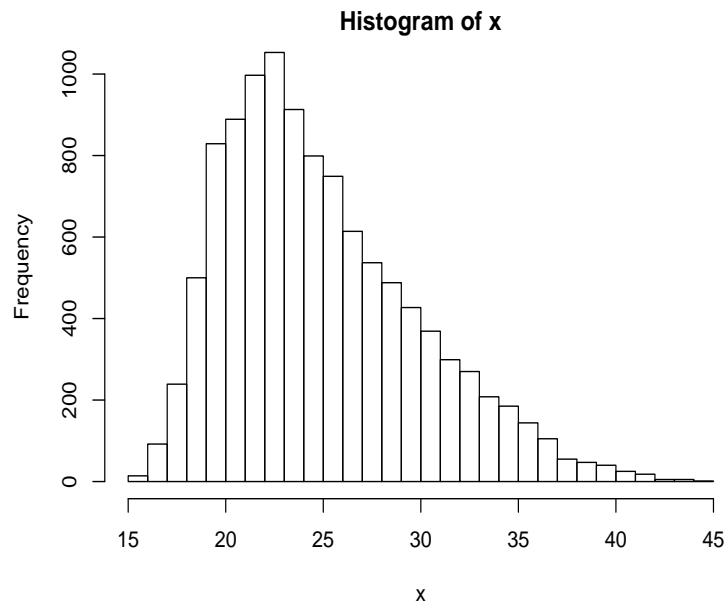
- ▶ proto je

$$\sigma_{\bar{X}}^2 = \sigma_{\frac{1}{n} \sum_{i=1}^n X_i}^2 = \frac{1}{n^2} n \sigma^2 = \frac{\sigma^2}{n}$$

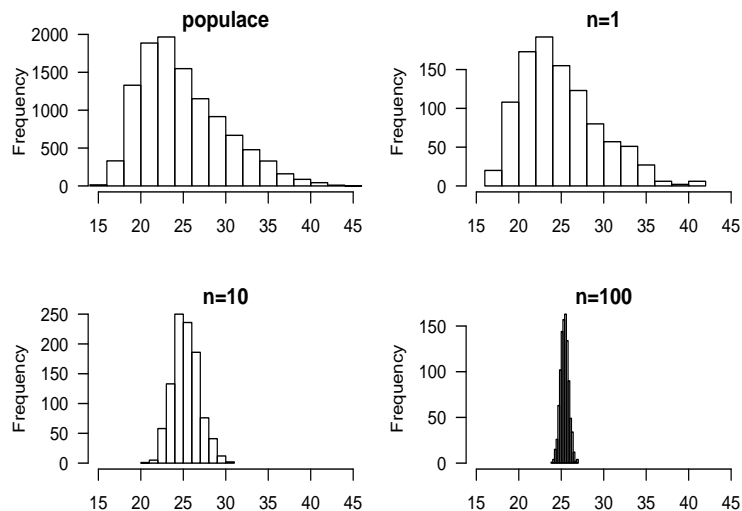
- ▶ průměr  $\bar{X}$  má tedy rozptyl  $n$ -krát menší, než jednotlivá pozorování
- ▶ **střední chyba** průměru = směrodatná odchylka průměru

$$\text{S.E.}(\bar{X}) = \frac{\sigma}{\sqrt{n}}$$

## příklad: věk matek



## příklad: histogram populace a histogramy výběrů šířky intervalů stejné

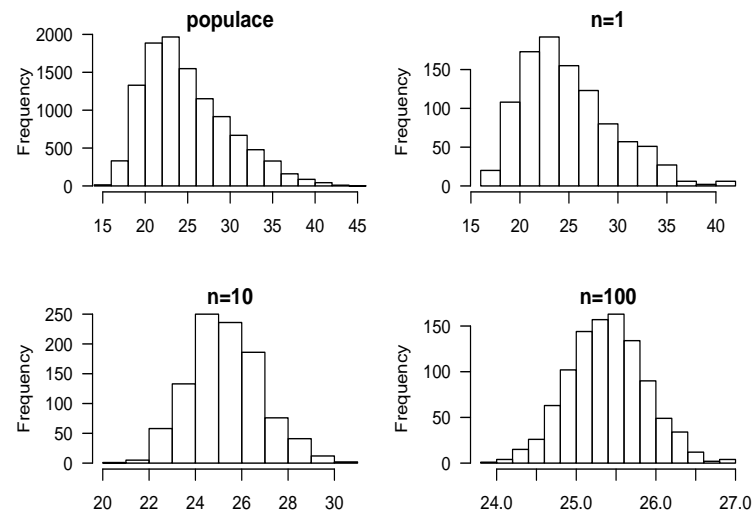


## příklad: věk matek

- ▶ výjimečný umělý příklad, kdy známe celou populaci
- ▶ populace obsahuje 10 916 hodnot
- ▶ rozdělení věku je výrazně nesymetrické
- ▶ prováděn výběr rozsahu  $n$ , vždy spočítán průměr
- ▶  $N$ krát opakovaně provedeno (spočítáno  $N = 1000$  průměrů)
- ▶ spočítány charakteristiky z  $N$  průměrů jako výchozích hodnot, (modře charakteristiky celé populace nebo hodnoty odvozené)

$n$	průměr	sm. odch.	$\sigma/\sqrt{n}$	šikmost	špičatost
1	25.43	4.62	4.94	0.74	0.29
10	25.35	1.54	1.56	0.28	-0.04
100	25.39	0.48	0.49	0.08	-0.05
(populace)	$\mu = 25.40$	$\sigma = 4.94$	4.94	0.77	0.19

## příklad: histogram populace a histogramy výběrů šířky intervalů přizpůsobené variabilitě



## příklad: shrnutí

- ▶ průměry kolísají kolem populačního průměru  $\mu$
- ▶ směrodatné odchylky klesají s rostoucím  $\sqrt{n}$
- ▶ šikmost a špičatost se s rostoucím  $n$  blíží k nule
- ▶ je naděje, že s rostoucím  $n$  je histogram podobnější hustotě normálního rozdělení – projev *centrální limitní věty*

interval spolehlivosti pro populační průměr  $\mu$ 

- ▶ pro nezávislé náhodné veličiny  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$  platí

$$\bar{X} \sim N(\mu, \sigma^2/n)$$

- ▶ proto je  $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$

- ▶ použijeme kritickou hodnotu

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

- ▶ hodnota parametru  $\mu$  je tedy s pstí  $1 - \alpha$  pokryta intervalem

$$\left(\bar{X} - \frac{\sigma}{\sqrt{n}}z(\alpha/2); \bar{X} + \frac{\sigma}{\sqrt{n}}z(\alpha/2)\right)$$

- ▶ lze použít pro velká  $n$  i bez požadavku na normální rozdělení

## centrální limitní věta

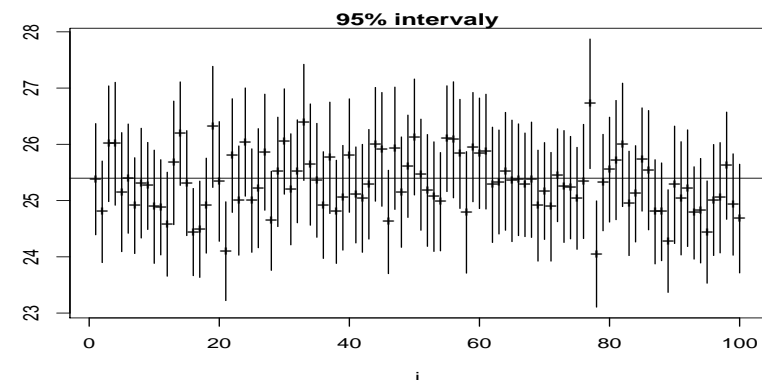
- ▶ vlastnost součtu nezávislých náhodných veličin se stejným rozdělením (populační průměr  $\mu$ , popul. rozptyl  $\sigma^2$ )
- ▶ průměr je součet dělený počtem sčítanců  
⇒ pro průměr platí CLV také
- ▶ standardizovaný součet (průměr)  $n$  nezávislých náhodných veličin lze pro velké  $n$  aproximovat normálním rozdělením  $N(0, 1)$

$$Z = \frac{\sum_{i=1}^n X_i - n \cdot \mu}{\sigma \sqrt{n}} = \frac{\bar{X} - \mu}{\sigma} \sqrt{n} \sim N(0, 1)$$

▶ CLV pro četnosti

- ▶ pro velká  $n$  se výběrový průměr chová, jako by šlo o výběr z normálního rozdělení, a to bez ohledu na výchozí rozdělení

$$\bar{X} \sim N(\mu, \sigma^2/n)$$

100 intervalů spolehlivosti ( $n = 100$ ,  $1 - \alpha = 95\%$ )  
(v 7 případech interval **neobsahuje**  $\mu$ )

## příklad: IQ vysokoškoláků

- ▶ u  $n = 16$  náhodně vybraných studentů jisté fakulty byla zjištěna hodnota IQ
- ▶ metoda měření IQ je konstruována tak, že je  $\sigma = 15$
- ▶ vyšel průměr  $\bar{x} = 110$
- ▶ co lze říci o populačním průměru všech studentů oné velké fakulty?
- ▶ 95% interval spolehlivosti ( $z(0,025) = 1,96$ ):

$$\left(110 - \frac{15}{4} \cdot 1,96; 110 + \frac{15}{4} \cdot 1,96\right) = (102,65; 117,35)$$

- ▶ skutečný populační průměr  $\mu$  (všech studentů oné fakulty) leží s 95% pravděpodobností mezi 102,65 a 117,35
- ▶  $\mu$  leží s 90% pravděpodobností mezi 103,83 a 116,17

interval spolehlivosti pro  $\mu$  (neznámé  $\sigma$ )

- ▶ neznáme-li  $\sigma$ , nahradíme je pomocí (výběrová směr. odchylka)

$$S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}$$

- ▶ interval spolehlivosti pro  $\mu$ :

$$\left(\bar{X} - \frac{S}{\sqrt{n}} t_{n-1}(\alpha); \bar{X} + \frac{S}{\sqrt{n}} t_{n-1}(\alpha)\right)$$

- ▶ použití kritické hodnoty  $t_{n-1}(\alpha)$  Studentova  $t$ -rozdělení místo kritické hodnoty  $z(\alpha/2)$  je penalizací za to, že neznámou směrodatnou odchylku  $\sigma$  jsme nahradili jejím odhadem  $S$
- ▶ platí totiž  $t_{n-1}(\alpha) > z(\alpha/2)$ , s rostoucím  $n$  se rozdíl zmenšuje

vlastnosti intervalu spolehlivosti pro  $\mu$ 

- ▶ délka intervalu roste s požadovanou spolehlivostí
  - ▶ 90% interval (103,83; 116,17) má délku 12,34
  - ▶ 95% interval (102,65; 117,35) má délku 14,70
- ▶ délka intervalu klesá s rostoucím počtem pozorování  $n$ 
  - ▶ pro  $n = 16$  má 95% interval (102,65; 117,35) délku 14,70
  - ▶ pro  $n = 16 \cdot 4 = 64$  má 95% interval (106,325; 113,675) délku 7,35, tedy poloviční
- ▶ kolik potřebujeme pozorování, aby měl 95% interval délku  $2\delta$ ?

$$\frac{\sigma}{\sqrt{n}} z(\alpha/2) = \delta \quad \Rightarrow \quad n = \left(\frac{\sigma}{\delta} z(\alpha/2)\right)^2$$

- ▶ v příkladu s IQ požadujeme  $\delta = 1$ :

$$n = \left(\frac{15}{1} 1,96\right)^2 \doteq 864$$

## příklad: výška postavy

- ▶ studenti odhadovali výšku přednášejícího; předpokládejme, že nestranně a nezávisle na sobě
- ▶  $n = 22$ ,  $\bar{x} = 172,4$ ,  $s_x = 4,032$
- ▶ z tabulek:  $t_{21}(0,05) = 2,080$

$$\left(172,4 - \frac{4,032}{\sqrt{22}} \cdot 2,080; 172,4 + \frac{4,032}{\sqrt{22}} \cdot 2,080\right) \\ (170,6; 174,2)$$

- ▶ skutečná výška je s pravděpodobností 95 % někde mezi 170,7 cm a 174,2 cm
- ▶  $z(0,025) = 1,96$

## centrální limitní věta pro četnosti

- ▶ co říká CLV? CLV
- ▶ absolutní četnost  $Y$ 
  - ▶  $Y$  – součet nezávislých veličin s alternativním rozdělením
  - ▶ populační průměr  $X_i$  je  $\pi$
  - ▶ populační rozptyl  $X_i$  je  $\pi(1 - \pi)$
  - ▶  $Y = \sum_{i=1}^n X_i$
  - ▶  $Y \sim \text{bi}(n, \pi)$ , proto přibližně  $Y \sim N(n\pi, n\pi(1 - \pi))$
- ▶ relativní četnost  $f = Y/n$ 
  - ▶  $f$  – průměr nezávislých veličin s alternativním rozdělením
  - ▶  $f \sim N(\pi, \pi(1 - \pi)/n)$

interval spolehlivosti pro podíl (pravděpodobnost)  $\pi$ 

- ▶  $\pi$  – podíl prvků populace s danou vlastností
- ▶  $\pi$  – pst, s jakou takový prvek vylosujeme
- ▶ počet prvků náhodně vybraných s onou vlastností  $Y \sim \text{bi}(n, \pi)$
- ▶ střední chyba relativní četnosti  $Y/n = f$   
= směrodatná odchylka relativní četnosti  $f$   
= odmocnina z rozptylu relativní četnosti  $f$  je tedy  $\sqrt{\frac{\pi(1-\pi)}{n}}$
- ▶ pravděpodobnost  $\pi$  neznáme, odhadneme ji pomocí  $f$
- ▶ odtud je přibližný 95% interval spolehlivosti pro  $\pi$

$$\left( f - 1,96 \cdot \sqrt{\frac{f(1-f)}{n}}; f + 1,96 \cdot \sqrt{\frac{f(1-f)}{n}} \right)$$

- ▶ skutečná pst  $\pi$  je tedy s 95% pstí v uvedeném rozmezí
- ▶ existuje přesnější (pracnější) postup

## příklad: počet studentek

- ▶ za zkušenosti je známo, že mezi uchazeči o studium bývá 45 % dívek
  - ▶ s jakou pravděpodobností bude při 500 přihláškách počet dívek mezi 200 a 220 (včetně)?
  - ▶  $Y \sim \text{bi}(500, 0,45)$  má  $\mu_Y = 500 \cdot 0,45 = 225$ ,  
 $\sigma_Y^2 = 500 \cdot 0,45 \cdot 0,55 = 123,75$ , tedy  $\sigma_Y = 11,1$
- $$P(200 \leq Y \leq 220) \doteq \Phi\left(\frac{220,5 - 225}{11,1}\right) - \Phi\left(\frac{199,5 - 225}{11,1}\right)$$
- ▶ hledaná pravděpodobnost je přibližně 33,2 % (přesně 33,3 %)  
[NORMDIST(220,5;225;11,1243;1)  
-NORMDIST(199,5;225;11,1243;1)]  
[pnorm(220.5,500\*0.45,sqrt(500\*0.45\*0.55))  
-pnorm(199.5,500\*0.45,sqrt(500\*0.45\*0.55))]  
[BINOMDIST(220;500;0,45;1)-BINOMDIST(199;500;0,45;1)]  
[pbinom(220,500,0.45)-pbinom(199,500,0.45)]

## příklad: hody s hrací kostkou

- ▶ odhadujeme pravděpodobnost šestky
- ▶ kostka A:  $n = 100$ ,  $n_A = 17$ ,  $f_A = 0,17$

$$\left( 0,17 - 1,96 \cdot \sqrt{\frac{0,17 \cdot 0,83}{100}}; 0,17 + 1,96 \cdot \sqrt{\frac{0,17 \cdot 0,83}{100}} \right)$$

**(0,10; 0,24)**

- ▶ kostka B:  $n = 100$ ,  $n_B = 41$ ,  $f_B = 0,41$

$$\left( 0,41 - 1,96 \cdot \sqrt{\frac{0,41 \cdot 0,59}{100}}; 0,41 + 1,96 \cdot \sqrt{\frac{0,41 \cdot 0,59}{100}} \right)$$

**(0,31; 0,51)**

- ▶ důležitý rozdíl: u kostky A patří  $1/6 = 0,167$  do intervalu spolehlivosti; u kostky B nikoliv; může to něco znamenat?