

PSY117

Statistická analýza dat v psychologii

**Přednáška 8 2016**

---

# Statistické usuzování, odhady

Věci, které můžeme přímo pozorovat, jsou téměř vždy pouze vzorky.

*Alfred North Whitehead*

# Barevná srdíčka kolegyně Michalčákové

---

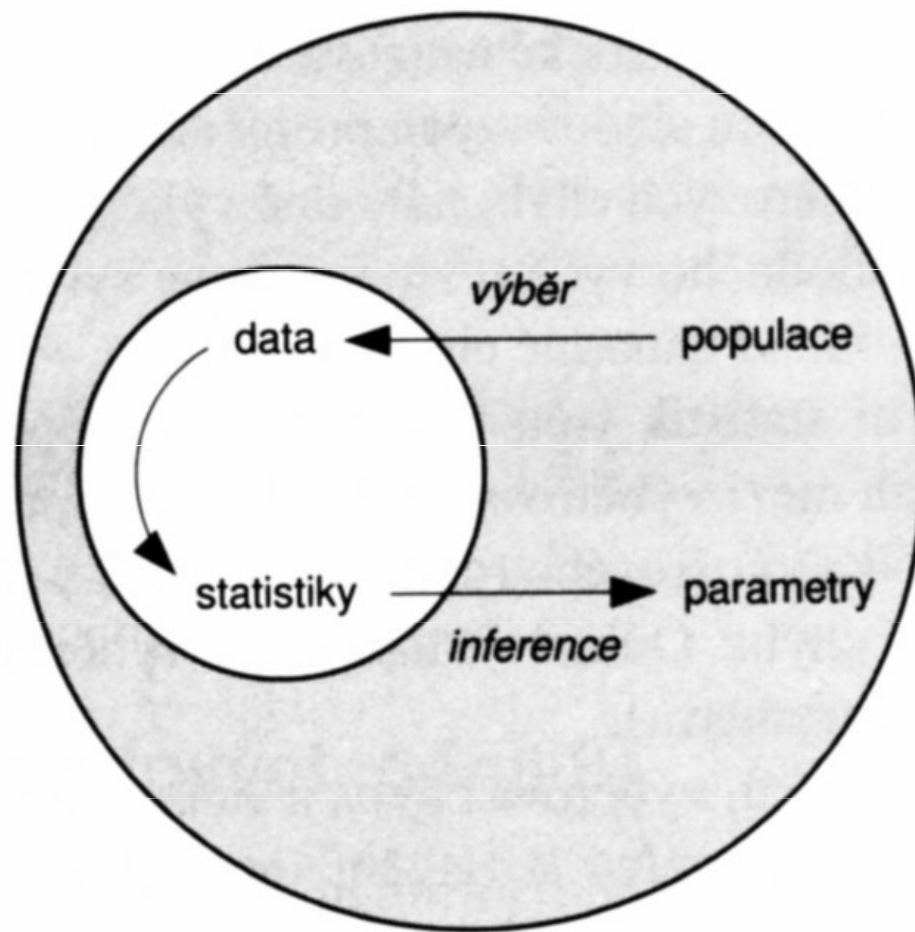
□ Jaký je podíl bílých a barevných srdíček v balení?

□ Simulace binomického rozložení

---

# Výběr – od deskripce k indukci

---



- Deskripce dat, odhad parametrů
- Usuzování = inference = indukce
- Počítá se s **náhodným výběrem**
  - tj. výběr jedince splňuje podmínky náhodného pokusu
  - není-li výběr v pravém slova smyslu náhodný, uvažujeme, v čem se p-dobně liší od náhodného

# Statistiky a parametry

---

- Na vzorku (datech) počítáme **statistiky**
- Hodnotě statistiky v celé populaci říkáme **parametr**.
  - Pro parametry používáme odpovídající písmena řecké abecedy
    - např. průměr: statistika  $m$ , parametr  $\mu$  (mí)
    - další:  $s - \sigma$  (sigma),  $r - \rho$  (ró),  $d - \delta$  (delta - rozdíl)
- Statistiky jsou **odhady** parametrů
  - tj. jsou vždy zatíženy chybou – **výběrovou chybou**
  - *chyby náhodné* – umíme spočítat, známe-li **výběrové rozložení**
  - *chyby systematické* – nevhodné statistiky, špatné měření, špatný způsob výběru vzorku (metodologie)

Jak dobré jsou tyto odhady?

# Výběrové rozložení a sm. chyba

---

- Spočítáme-li tutéž statistiku na mnoha nezávislých náhodných vzorcích
  - získáme mnoho různých odhadů parametru
  - tyto odhady mají nějaké rozložení - **výběrové rozložení (statistiky)**

[http://onlinestatbook.com/stat\\_sim/sampling\\_dist/index.html](http://onlinestatbook.com/stat_sim/sampling_dist/index.html)

- **Výběrové rozložení** statistik obvykle můžeme popsat
  - průměrem – ten se u dobrých statistik blíží hodnotě **parametru**
  - směrodatnou odchylkou – říkáme jí **směrodatná chyba** ((odhadu parametru) nebo také střední chyba a obecněji i výběrová chyba
  - Čím je velikost vzorku/ů větší, tím je směrodatná chyba menší

# Výběrové rozložení (odhadu) průměru

---

Odhad průměru má přibližně **normální rozložení**,

- jehož průměr je  $\mu$  se směrodatnou chybou .....  $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{N}}$
- Platí to i tehdy, když rozložení proměnné není normální.
  - a to „díky“ **centrálnímu limitnímu teorému**
- Jenomže my obvykle neznáme  $\sigma$ ...

Neznáme-li  $\sigma$ , musíme použít  $s$

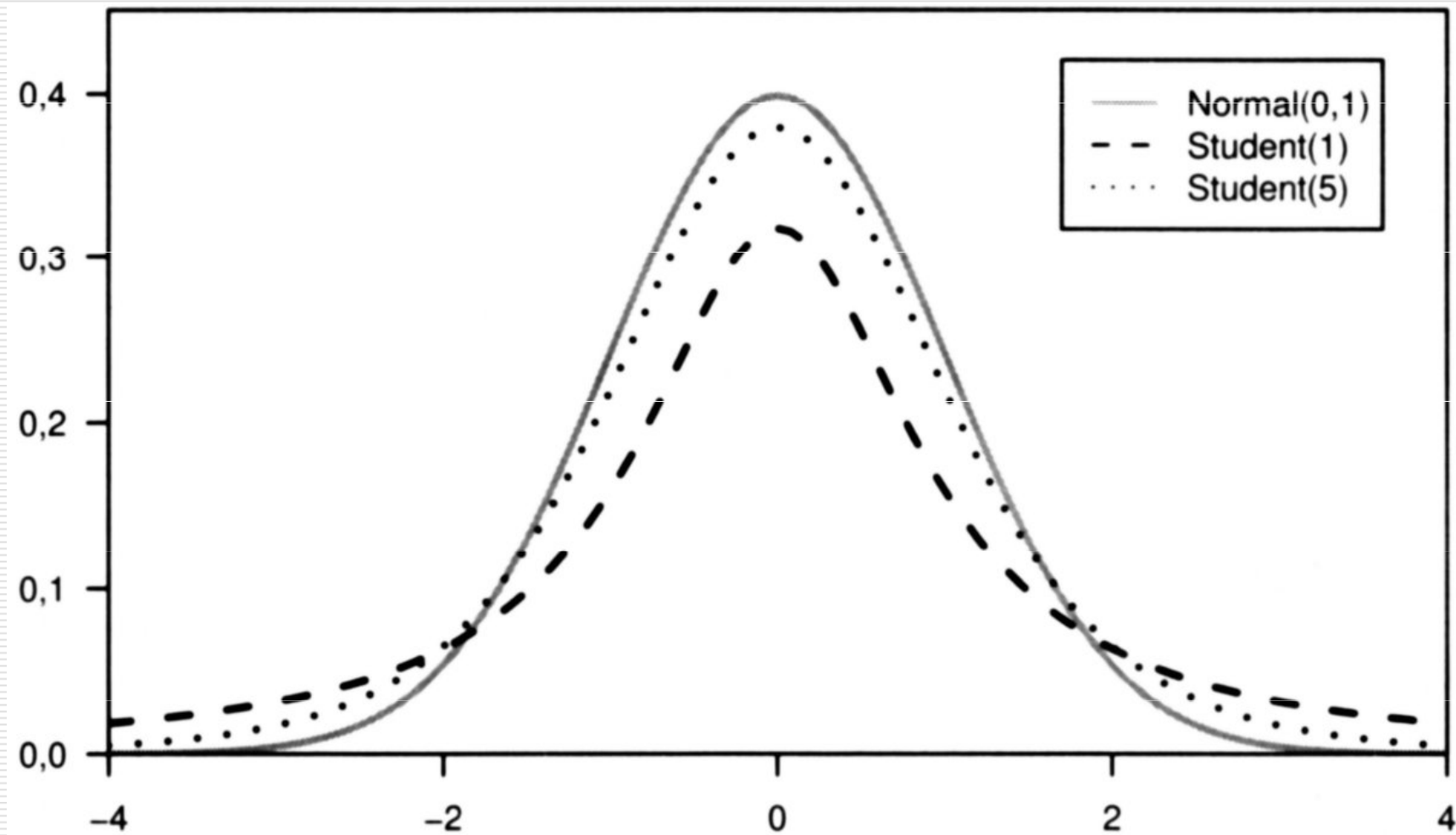
- průměr zůstává  $\mu$ , směrodatná chyba je nyní .....  $s_{\bar{x}} = \frac{s}{\sqrt{N}}$
- výběrové rozložení není normální, jde o

## **Studentovo $t$ -rozložení**

- jako normální s těžšími konci ( $t$  je pro  $t$ -rozložení totéž, co  $z$  pro normální rozložení)
- má různé tvary pro různá  $n$  : stupně volnosti –  $\nu$  (ný)
  - zde  $\nu = N-1$ ; čím vyšší  $N$ , tím se  $t$ -rozložení blíží normálnímu

# Studentovo $t$ -rozložení

---



# Výběrová rozložení dalších statistik

---

Nyní je tedy třeba ke každé popisné statistice znát ještě další vlastnost – její teoretické **výběrové rozložení**

- relativní četnost – přibližně normální - Hendl 162
- rozptyl – po transformaci  $\chi^2$ -rozložení (chí kvadrát) - Hendl 159
- Pearsonova  $r$  – po Fisherově transformaci normální – Hendl 252

Teoretická výběrová rozložení různých statistik jsou různá

- Statistika je obvykle transformována do podoby, která má jedno z běžných teoretických rozložení: normální, chí-kvadrát rozložení (Pearsonovo),  $t$ -rozložení (Studentovo),  $F$ -rozložení (Fisherovo, Snedecorovo)
- Netřeba je znát z hlavy, programy je používají za vás, ale stojí za to vědět, že existují přehledy – např. Receptář Oseckých nebo Sheskin ISBN 1584884401
- Pro interpretační potřeby si obvykle vystačíme s představou výběrového rozložení průměru
- Pozor, centrální limitní teorém se týká pouze výběrového rozložení průměru!



# Estimační kvality statistik I

Kvality statistiky jako prostředku odhadu „skutečné“ hodnoty v populaci

**TABLE 5.1**

The Expected Values of the Range,  $s^2$ , and  $s$  as a Function of Sample Size of  $n$  Observations from a Random Sample from a Normal Distribution in which  $\sigma = 10$

<i>If <math>\sigma = 10</math> <math>n</math></i>	<i>Expected Value of the Range</i>	<i>Expected Value of <math>s^2</math></i>	<i>Expected Value of <math>s</math></i>	<i>Expected Value of Range/<math>s</math></i>
2	11	100	8.0	1.4
5	23	100	9.4	2.4
10	31	100	9.73	3.2
20	37	100	9.87	3.7
50	45	100	9.95	4.5
100	50	100	9.97	5.0
200	55	100	9.987	5.5
500	61	100	9.993	6.1
1,000	65	100	9.997	6.5

# Estimační kvality statistik II

---

- Nezkreslenost
  - tj. že systematicky nenad(pod)hodnocuje
  - např. s podhodnocuje
- Konzistence
  - s velikostí vzorku roste přesnost odhadu
- Relativní účinnost
  - jak rychle roste přesnost s velikostí vzorku
  - zde vítězí  $M$  nad  $Md$  a strhává s sebou i další momentové statistiky
    - jejich výhodou je i snadné počítání s nimi

*Alternativně Kvalita bodového odhadu viz Hendl 175*

---

# Bodové vs. intervalové odhady

$\alpha$  je p-nost chyby a proto je hladina spolehlivosti  $1-\alpha$ , tj. 95% spolehlivost znamená 5% chybovost:  $(1-0,05)$

Parametr se můžeme snažit odhadnout...

- **bodovým odhadem** – tj. odhadujeme přímo hodnotu parametru, např. průměr.
- **intervalovým odhadem** – tj. odhadnutím intervalu, který parametr s určitou p-ností zahrnuje
  - výsledkem intervalového odhadu je **interval spolehlivosti**
  - interval spolehlivosti tvoříme z bodového odhadu a znalosti jeho výběrového rozložení, tj. (bod $\pm$ odchylka)
  - intervalový odhad lepší - více informací  $(1-\alpha) CI = \bar{X} \pm z_{1-\alpha/2} \sigma_{\bar{X}}$
  - té p-nosti se v tomto kontextu říká **hladina spolehlivosti**  $(1-\alpha)$ 
    - typicky se používá 95% a 99% hladina spolehlivosti
    - pak říkáme, že hledaný parametr je s 95% p-ností v intervalu spolehlivosti

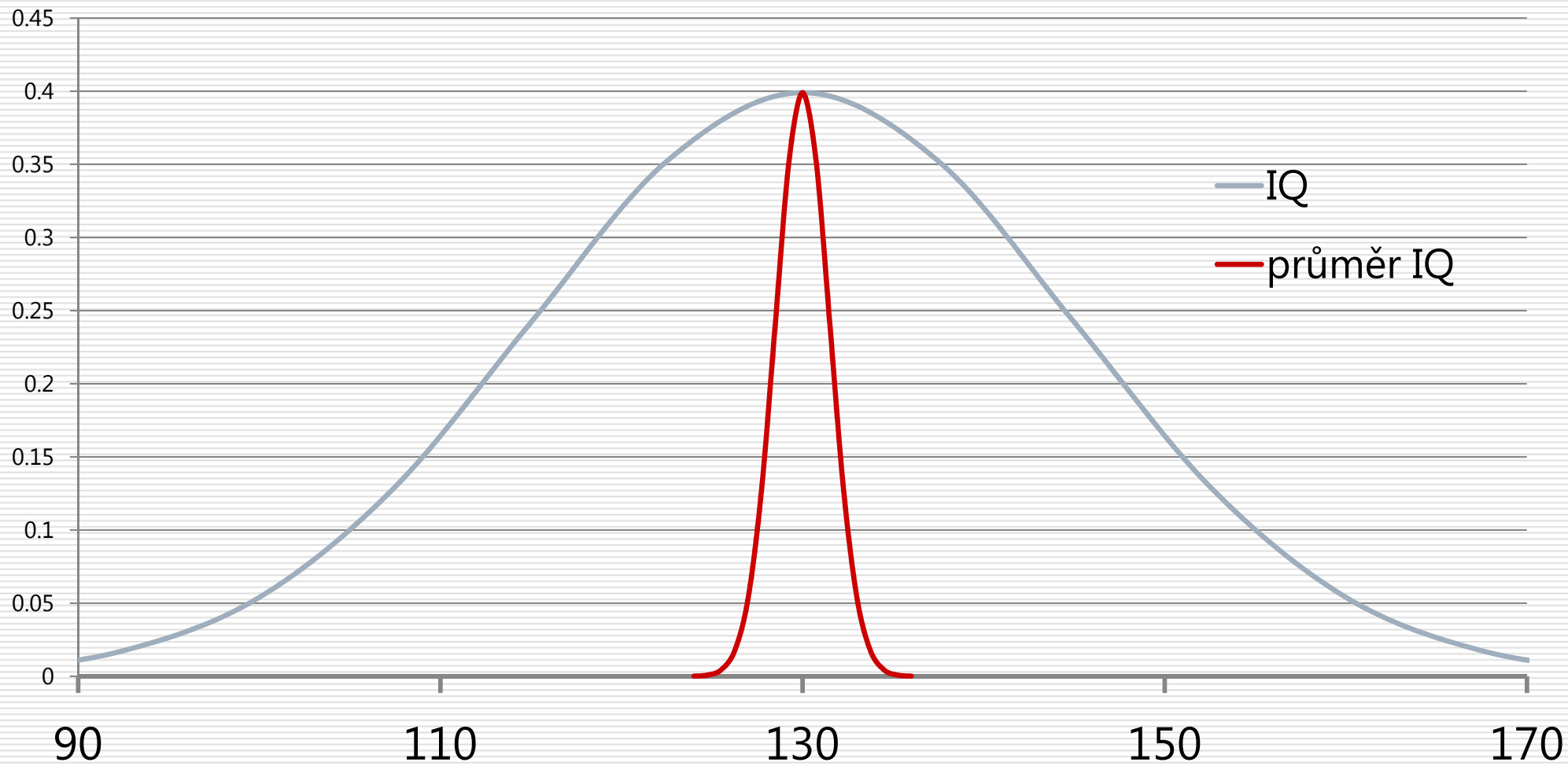
Zkuste si sami: [http://onlinestatbook.com/stat\\_sim/conf\\_interval/index.html](http://onlinestatbook.com/stat_sim/conf_interval/index.html)

# Příklad konstrukce intervalu spolehlivosti pro průměr 1

---

**Na vzorku dětí ( $N=100$ ) s různobarevnými očima jsme spočítali průměrné IQ 130, přičemž víme, že  $\sigma = 15$ .**

- **bodový odhad** průměrného IQ v populaci dětí s různobarevnými očima (tj. parametru,  $\mu$ ) je 130
  - **intervalový odhad**
    - Známe-li  $\sigma$ , výběrové rozložení průměru má **normální rozložení...**
    - ...se středem v  $\mu$ .  $\mu$  neznáme, a tak použijeme bodový odhad  $m = 130$
    - ... se směrodatnou chybou odhadu průměru  $s_m = \sigma/\sqrt{N} = 15/\sqrt{100} = 1,5$ .
    - Zvolíme-li hladinu spolehlivosti  $1-\alpha = 95\%$ ,
    - pak v tabulkách/Excelu zjistíme, že 95% normálního rozl. je mezi hodnotami  $z = -1,96$  a  $1,96$ , tj.  $1-\alpha/2 z = 0,975 z = 1,96$ , Excel: =NORMSINV(0,975)
    - interval spolehlivosti:  $(m - 1,96s_m; m + 1,96s_m) = (127,1; 132,9)$ ,
    - **tj. s 95% pravděpodobností  $127,1 \leq \mu \leq 132,9$**
-



# Příklad konstrukce intervalu spolehlivosti pro průměr 2

---

**Na vzorku dětí ( $N=100$ ) s různobarevnými očima jsme spočítali průměrné IQ 130 a  $s = 15$ .**

- **bodový odhad** průměrného IQ v populaci dětí s různobarevnými očima (tj. parametru,  $\mu$ ) je 130
- **intervalový odhad**
  - střed intervalu spolehlivosti bude na bodovém odhadu, tj.  $m = 130$
  - víme, že výběrové rozložení průměru má  $t$ -rozložení se stupni volnosti  $\nu = N - 1 = 99$
  - zvolíme-li hladinu spolehlivosti  $1 - \alpha = 95\%$ ,
  - pak v tabulkách (Excelu) zjistíme, že 95%  $t$ -rozložení je mezi hodnotami  $t = -1,98$  a  $1,98$  (tj.  $_{1-\alpha/2}t(\nu) = {}_{0,975}t(99) = 1,98$  excel: `TINV(0,05,99)`)
  - směrodatná chyba odhadu průměru  $s_m = s / \sqrt{n} = 15 / \sqrt{100} = 1,5$
  - interval spolehlivosti:  $(m - 1,98s_m; m + 1,98s_m) = (127,0 ; 133,0)$ ,
  - **tj. s 95% pravděpodobností  $127,0 \leq \mu \leq 133,0$**

pozor na tento rozdíl: ve středu intervalu je  $m$ , někde v intervalu je v 95% případech  $\mu$

# Interpretace intervalu spolehlivosti

---

- ... je prostá, avšak zrádná
  - 95% interval spolehlivosti znamená, že sestrojíme-li tento interval dle výše uvedených instrukcí, **v 95% případů sestojení intervalu tento interval zahrnuje odhadovaný parametr**, tj. v 95% případů je závěr, že  $\mu$  je mezi čísly  $a$  a  $b$ , správný.
  - V tomto smyslu to také znamená, že máme subjektivní 95% jistotu, že parametr je v námi určeném intervalu.
  - V konkrétním případě, kdy jsme spočetli konkrétní interval spolehlivosti ( $127 \leq \mu \leq 133$ ), to neznamená, že v 95% případech je  $\mu$  v intervalu od 127 do 133.
    - To proto, že  $\mu$  je konstanta; při opakovaných výzkumech se nemění. Díky omylnému výběru v každém výzkumu vychází poněkud jiný interval sestojený podle jiného výběrového průměru. Jinými slovy, trefujeme se obručí na kolík a ne kolíkem do obruče.
  - O čem tohle slovíčkaření je? O rozdílu mezi četnostním a subjektivním (Bayesovským) pojetím pravděpodobnosti.
-

# ...Výběrové rozložení mediánu

---

- Simulace: [www.stat.tamu.edu/~jhardin/applets/signed/SampDist2.html](http://www.stat.tamu.edu/~jhardin/applets/signed/SampDist2.html)
- V případě normálního rozložení je taky normální a směrodatná chyba je cca 1,25 směrodatné chyby průměru
- Pořadový způsob nabízí Campbell a Gardner<sup>1</sup>
  - Přibližný interval (pro  $N > 100$ ) se stanovuje opravdu pořadovým způsobem, tj. počítáme pořadí, které určuje horní a dolní mez intervalu
  - Pro 95% interval spolehlivosti pak je  $r$  pořadí určující horní mez a  $s$  pořadí určující dolní mez

$$r = \frac{n}{2} - z_{1-\alpha/2} \frac{\sqrt{n}}{2} \qquad s = 1 + \frac{n}{2} + z_{1-\alpha/2} \frac{\sqrt{n}}{2}$$

- Bootstrap
  - Obecná metoda, nejen pro mediány, téměř bez předpokladů (neparametrická)
  - Algoritmus:
    - 1. Proveďte výběr s navrácením ze svého výběru (o velikosti  $N$ )
    - 2. Spočítejte medián a uložte
    - 3. Opakujte kroky 1 a 2 tisíckrát
  - 95% interval je ohraničen 25. a 975. nejvyšším spočítaným mediánem.

---

<sup>1</sup>Campbell, M.J., Gardner, M.J. (2000). Medians and their differences. In Altman et al., *Statistics with confidence* (36 – 44). BMJ Books.



## ...Výběrové rozložení **relativní četnosti** $p$

---

- Pro dostatečně velkou populaci ( $np > 10$ ;  $n(1-p) > 10$ )...
- ...je přibližně normální s průměrem  $p$  a směrodatnou chybou  $\sqrt{p(1-p)/n}$
- $(1-\alpha)\%$  interval spolehlivosti má tedy podobu:

$$\left( p - z_{1-\alpha/2} \sqrt{p(1-p)/n}; p + z_{1-\alpha/2} \sqrt{p(1-p)/n} \right)$$

---

## ...Výběrové rozložení **rozptylu** $s^2$

---

- Rozložení poměru  $(s^2/\sigma^2)(n-1)$  má podobu chí-kvadrát rozložení s  $\nu = n-1$  stupni volnosti

$$\frac{s^2}{\sigma^2} (n - 1) \sim \chi^2(\nu)$$

- $(1-\alpha)\%$  interval spolehlivosti pro  $\sigma^2$  má tedy podobu:

$$\left( s^2 \frac{n-1}{\chi^2_{1-\alpha/2}(\nu)} ; s^2 \frac{n-1}{\chi^2_{\alpha/2}(\nu)} \right)$$

- V Excelu  $=\text{CHISQ.INV}(1-\alpha;df) = \chi^2_{1-\alpha}(df)$  [ $=\text{CHIINV}(\alpha;df)$ ]
-

## ...Výběrové rozložení Pearsonovy **korelace** $r$

---

- Výběrové rozložení korelace neznáme.
  - Známe výběrové rozložení korelace po Fisherově transformaci:  
 $Z = 0,5 \ln((1+r)/(1-r)) = \operatorname{arctgh}(r) = \operatorname{FISHER}(r)$
  - Výběrové rozložení  $Z$  je přibližně normální s průměrem  $Z$  a směrodatnou chybou  $s_Z = 1/\sqrt{n-3}$
  - $(1-\alpha)\%$  CI pro  $Z$ :  $(Z - z_{1-\alpha/2}s_Z; Z + z_{1-\alpha/2}s_Z)$
  - Nutno transformovat zpět do metriky korelačního koeficientu:  $r = (e^{2Z} - 1)/(e^{2Z} + 1) = \operatorname{FISHERINV}(Z)$   
 $(\operatorname{FISHERINV}(Z - z_{1-\alpha/2}s_Z); \operatorname{FISHERINV}(Z + z_{1-\alpha/2}s_Z))$
-

# Shrnutí

---

- Na vzorcích počítáme **statistiky**, které jsou odhadem populačních **parametrů**.
  - K posouzení přesnosti takového odhadu musíme znát **výběrové rozložení** statistiky, kterou k odhadu používáme, zejména jeho variabilitu – **směrodatnou chybu**.
  - Směrodatná chyba klesá především s velikostí vzorku a s variabilitou jevu v populaci.
  - Přesnost odhadu parametru sdělujeme prostřednictvím **intervalu spolehlivosti**.
-