

# Velmi jemný úvod do biomedicínské statistiky

B00364 Zdravotnická informatika

Lubomír Štěpánek<sup>1, 2</sup>



<sup>1</sup>Oddělení biomedicínské statistiky  
Ústav biofyziky a informatiky  
1. lékařská fakulta  
Univerzita Karlova v Praze



<sup>2</sup>Katedra biomedicínské informatiky  
Fakulta biomedicínského inženýrství  
České vysoké učení technické v Praze

28. května 2020

(2020) Lubomír Štěpánek, CC BY-NC-ND 3.0 (CZ)



Dílo lze dále svobodně šířit, ovšem s uvedením původního autora a s uvedením původní licence. Dílo není možné šířit komerčně ani s ním jakkoliv jinak nakládat pro účely komerčního zisku. Dílo nesmí být jakkoliv upravováno. Autor neručí za správnost informací uvedených kdekoli v předložené práci, přesto vynaložil nezanedbatelné úsilí, aby byla uvedená fakta správná a aktuální, a práci sepsal podle svého nejlepšího vědomí a svých „nejlepších“ znalostí problematiky.

# Obsah

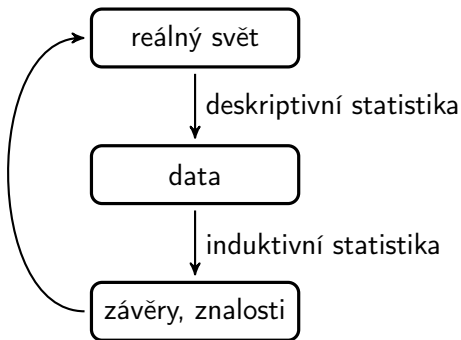
- 1 Úvod
- 2 Základní pojmy
- 3 Deskriptivní statistika
- 4 Pravděpodobnost
- 5 Induktivní statistika
- 6 Literatura



# Dělení statistiky

- deskriptivní statistika
  - popisuje data, ale nedělá na nich žádné „velké“ závěry
- induktivní statistika
  - pozoruje konkrétní data a vyvozuje z nich obecné závěry, ovšem s udáním stupně jejich spolehlivosti

# Vzájemný vztah deskriptivní a induktivní statistiky



# Pojem *statistický znak, veličina*

- statistický znak, veličina
  - měřitelná (veličina) či jinak zjistitelná (znak) charakteristika našeho zájmu
  - např. tělesná výška, pohlaví, mzda, apod.

# Pojem *statistická jednotka*

- statistická jednotka
  - základní atomický prvek zájmu, u nějž lze měřit nebo jinak získat hodnotu statistického znaku či veličiny
  - např. student, pacient, stát, molekula, apod.



# Pojem *statistický soubor*

- statistický soubor
  - množina statistických jednotek (prvků statistického souboru)
  - např. třída žáků, kohorta pacientů, apod.

## Vztah statistického znaku (veličiny), jednotky a souboru

- každá statistická jednotka (prvek) statistického souboru má svou hodnotu<sup>1</sup> určitého zkoumaného statistického znaku či veličiny (jde-li o měřitelný znak)
- např. *ve školní třídě změříme tělesnou výšku každého žáka*
  - *školní třída* je statistický soubor
  - *žáci* jsou statistické jednotky (prvky)
  - *tělesná výška* je statistická veličina

<sup>1</sup>ta může eventuálně chybět nebo být neznámá (missing value)

# Intermezzo

- měříme tělesné hmotnosti v kohortě pacientů-diabetiků na interním oddělení
- určíme, co je v takovém případě
  - statistickým znakem, resp. veličinou
  - statistickou jednotkou
  - statistickým souborem

# Cíle deskriptivní statistiky

- cílem je popsat soubor dat
  - číselně (resp. tabulkou)
  - graficky
- popisné číselné ukazatele i grafické přístupy se liší, pokud jde
  - o kvantitativní statistický znak (veličinu)
  - o kvalitativní statistický znak



## Dělení kvantitativního znaku (veličiny)

- dle spojitosti číselných hodnot
  - *spojitý* – hodnoty nabývají reálných čísel, nebo je na ně lze převést nějakou bijekcí
    - např. hmotnost, výška atd.
  - *diskrétní* – hodnoty jsou oddělená čísla obvykle ve smyslu počet či pořadí
    - např. počty pacientů atd.
- dle měřítka
  - *intervalová stupnice* – lze si smysluplně odpovědět, o kolik se dvě hodnoty liší, ale ne kolikrát
    - např. °C, datumy atd.
  - *poměrová stupnice* – lze si smysluplně odpovědět, o kolik se dvě hodnoty liší i kolikrát se liší
    - např. °K

---

- dle měřítka

- *nominální stupnice* – dvě či více vzájemně se vylučujících, rovnocenných tříd, které nelze uspořádat na číselné ose
  - např. pohlaví {muž, žena}
  - rodinný stav muže {svobodný, ženatý, rozvedený, vdovec, registrovaný}
- *ordinální stupnice* – kategorie je možné uspořádat vzestupně/sestupně, lze si smysluplně odpovědět, která hodnota je větší než jiná (ale ne o kolik, natož kolikrát)
  - např. pořadí v závodu, grade tumoru {1, 2, 3, 4} atd.







# Aritmetický průměr

- pro  $n$  čísel  $x_1, x_2, \dots, x_n$  spočítáme jejich aritmetický průměr jako

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n}$$

100

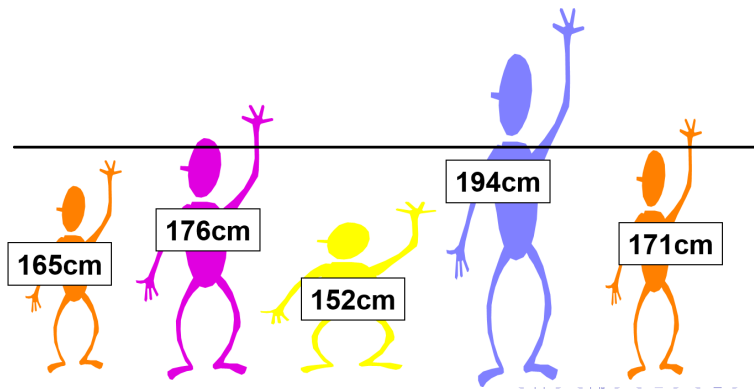






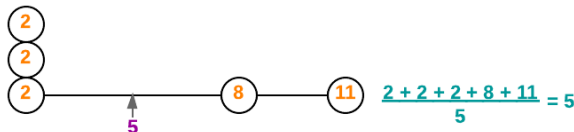
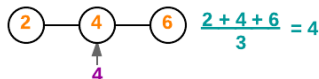
# Intermezzo

- určíme aritmetický průměr z následujícího souboru tělesných výšek
- $\bar{x} = \frac{165+176+152+194+171}{5} \doteq 171,6 \text{ [cm]}$
- kolik navzájem různých průměrů může mít jeden soubor čísel?
- pouze jeden



# Geometrická interpretace aritmetického průměru

- pokud zavěšíme  $n$  jednogramových závaží na pozice čísel  $x_1, x_2, \dots, x_n$  pravítka, hodnota průměru  $\bar{x}$  je v těžišti soustavy

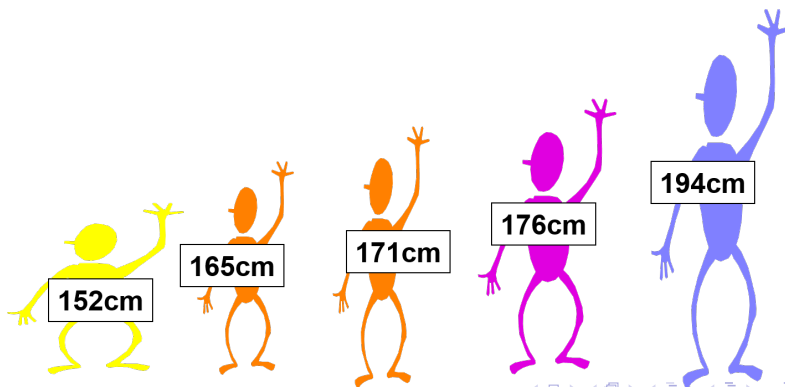






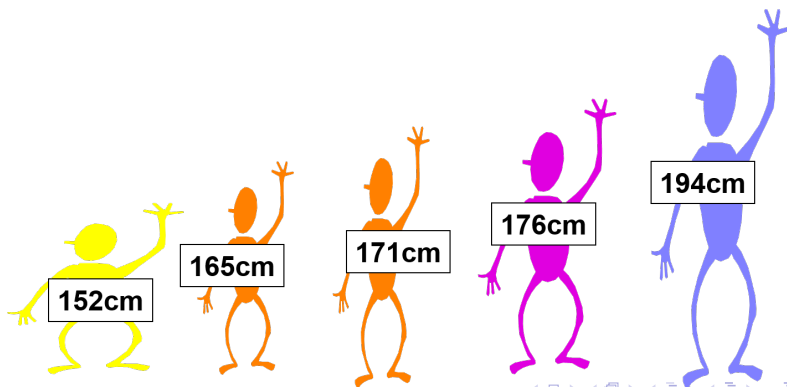
## Intermezzo

- určíme medián z následujícího souboru tělesných výšek



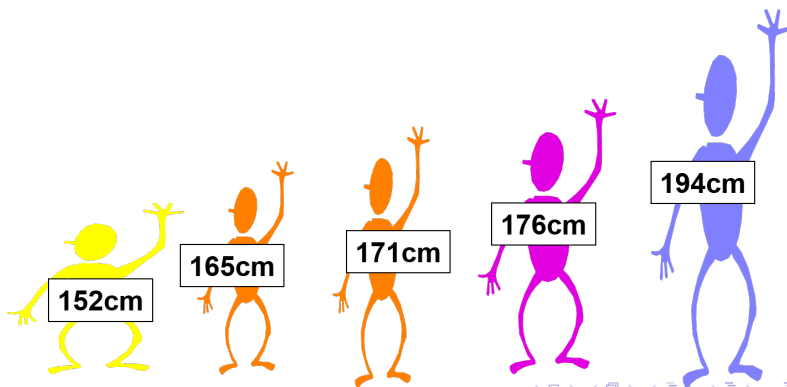
# Intermezzo

- určíme medián z následujícího souboru tělesných výšek
- $\tilde{x} = 171 \text{ [cm]}$



# Intermezzo

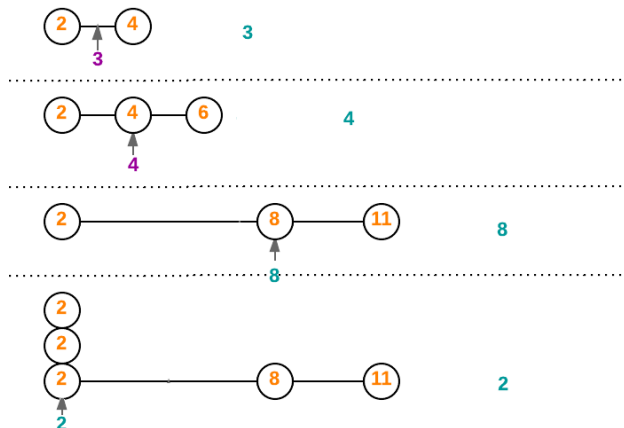
- určíme medián z následujícího souboru tělesných výšek
- $\tilde{x} = 171 \text{ [cm]}$
- kolik navzájem různých mediánů může mít jeden soubor čísel?





# Geometrická interpretace mediánu

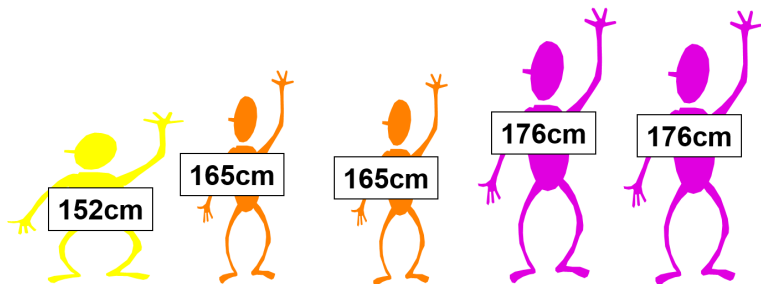
- pokud na pravítku vyznačíme pozice čísel  $x_1, x_2, \dots, x_n$ , hodnota mediánu  $\tilde{x}$  má od všech vyznačených bodů nejmenší možný součet vzdáleností





# Intermezzo

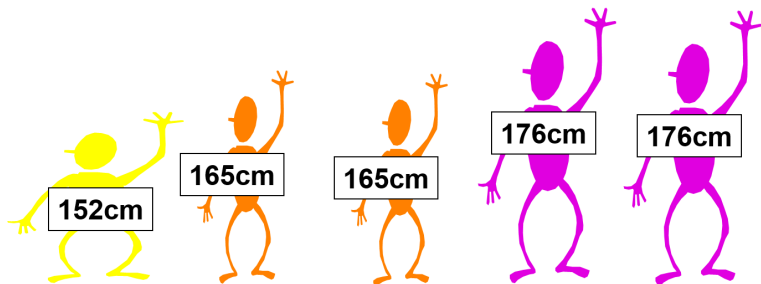
- určeme modus z následujícího souboru tělesných výšek





# Intermezzo

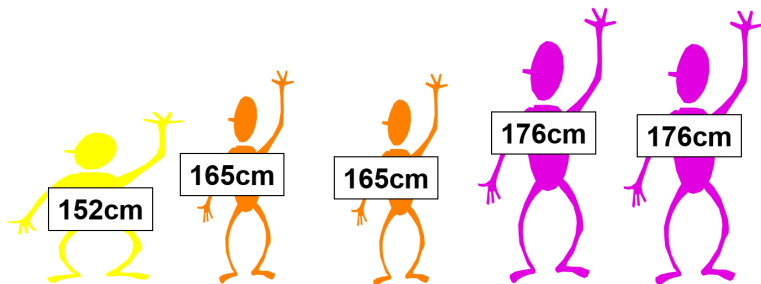
- určíme modus z následujícího souboru tělesných výšek
- $\hat{x} = \{165; 176\}$  [cm]



# Intermezzo

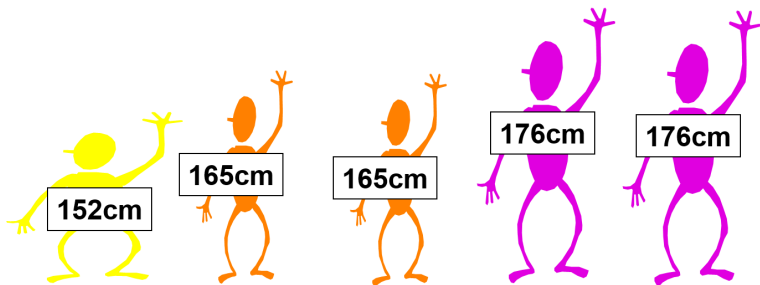
- určíme modus z následujícího souboru tělesných výšek
- $\hat{x} = \{165; 176\}$  [cm]
- kolik navzájem různých modů může mít jeden soubor čísel?

11



# Intermezzo

- určíme modus z následujícího souboru tělesných výšek
- $\hat{x} = \{165; 176\} \text{ [cm]}$
- kolik navzájem různých modů může mít jeden soubor čísel?
- alespoň jeden





# Intermezzo

- určíme aritmetický průměr a medián u každého z obou následujících souborů

$$\mathbf{x}_1 = \{1, 2, 3, 4, 5\} \quad \mathbf{x}_2 = \{1, 2, 3, 4, 90\}$$

- 

$$\bar{x}_1 = \tilde{x}_1 = 3; \quad \bar{x}_2 = 20; \tilde{x}_2 = 3$$

# Intermezzo

- určíme aritmetický průměr a medián u každého z obou následujících souborů

$$x_1 = \{1, 2, 3, 4, 5\} \quad x_2 = \{1, 2, 3, 4, 90\}$$

- 

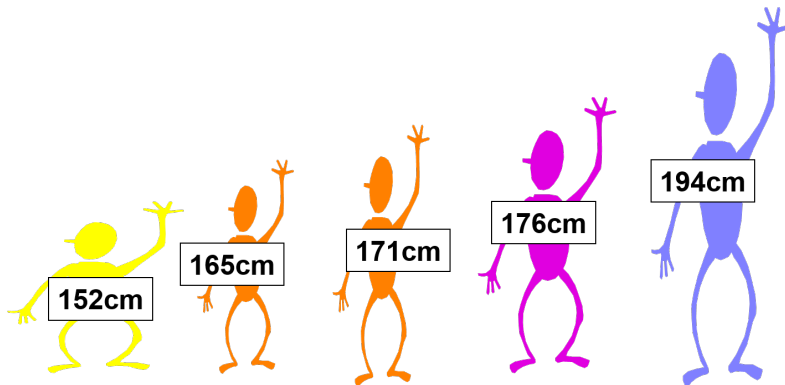
$$\bar{x}_1 = \tilde{x}_1 = 3; \quad \bar{x}_2 = 20; \quad \tilde{x}_2 = 3$$

- která z měr polohy (průměr, medián) lépe vyhovuje „asymetrickým“ datům?



# Intermezzo

- určeme rozpětí z následujícího souboru tělesných výšek



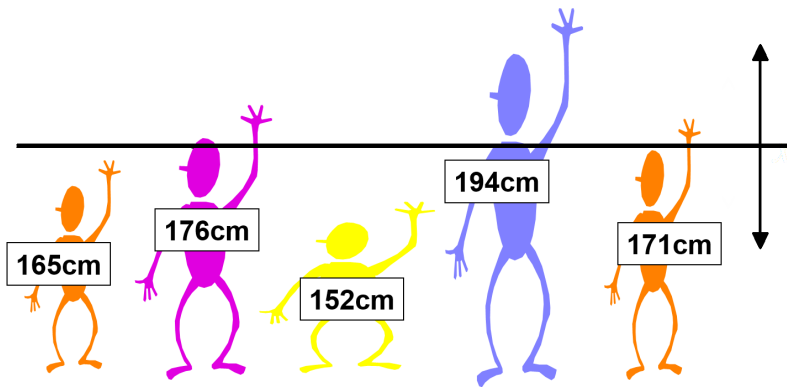






# Intermezzo

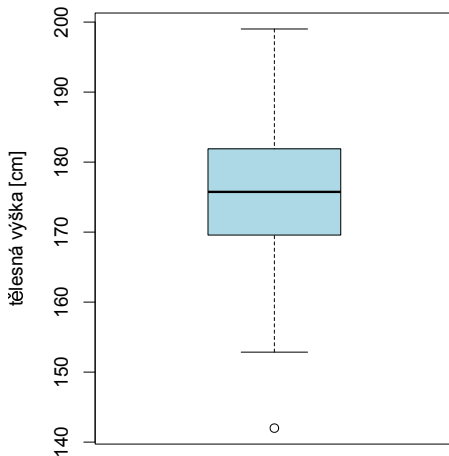
- určeme směrodatnou odchylku a rozptyl z následujícího souboru tělesných výšek





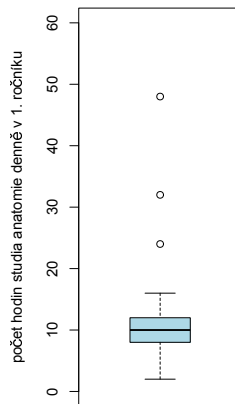
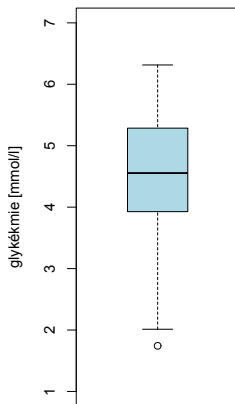
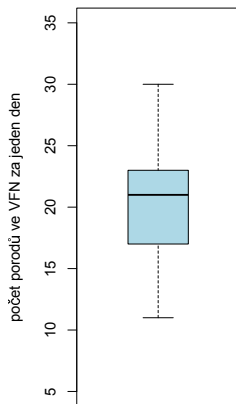
# Krabicový diagram (boxplot)

- vhodný pro kvantitativní znaky



## Intermezzo

- který z krabicových diagramů nedává smysl?



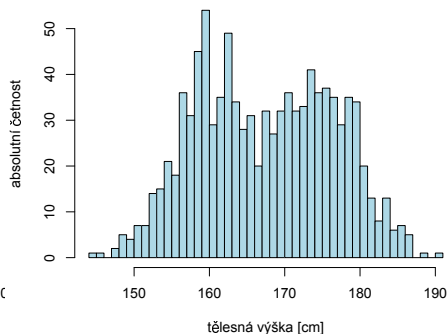
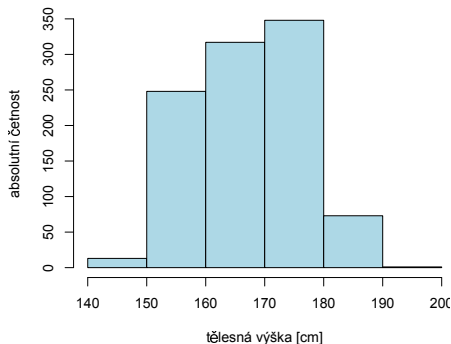


## Počet intervalů v histogramu

- rozdílný počet intervalů histogramu mění „příběh“ dat!
- nejčastěji je počet intervalů  $k$  dán Sturgesovým pravidlem

$$k = \lceil \log_2 n \rceil,$$

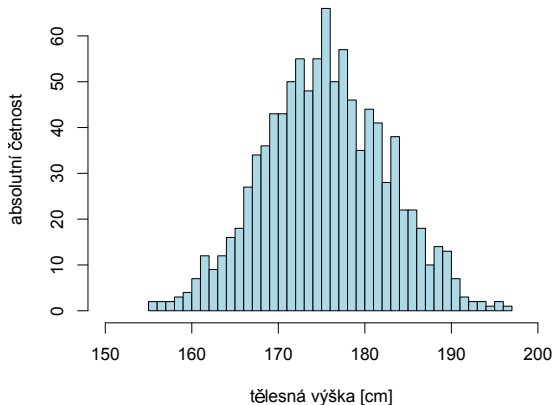
kde  $n$  je počet pozorování v souboru





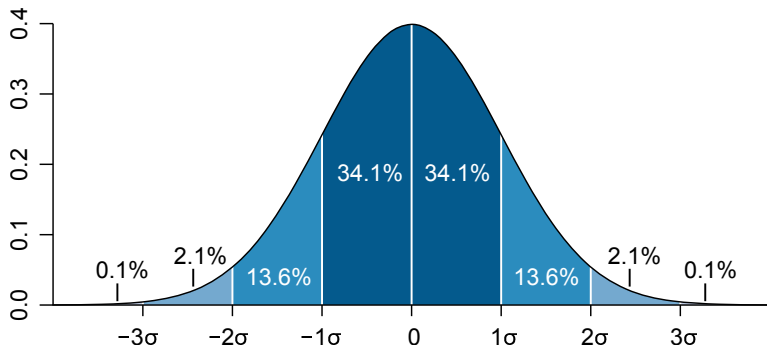
# Normální rozdělení kvantitativního znaku

- lze odhadnout z histogramu



## Vztah mezi polohou, variabilitou, tvarem a proporcí

- pokud je udržitelný předpoklad normálního rozložení, pak
  - v intervalu  $\langle \bar{x} - s, \bar{x} + s \rangle$  leží asi 68 % hodnot
  - v intervalu  $\langle \bar{x} - 2s, \bar{x} + 2s \rangle$  leží asi 95 % hodnot
  - v intervalu  $\langle \bar{x} - 3s, \bar{x} + 3s \rangle$  leží asi 99,7 % hodnot



---

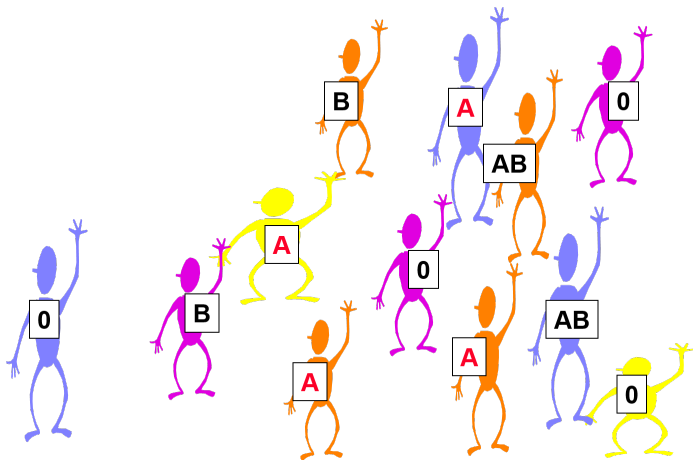
# Popis kvalitativního znaku

- např. krevní skupiny, grading tumoru, pohlaví, atd.
- číselně
  - absolutní, relativní četnosti
- graficky
  - koláčový diagram



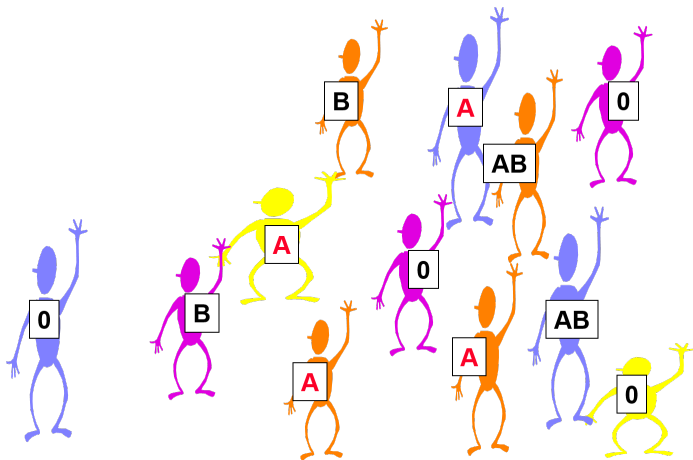
# Intermezzo

- určeme absolutní a relativní četnost krevní skupiny A



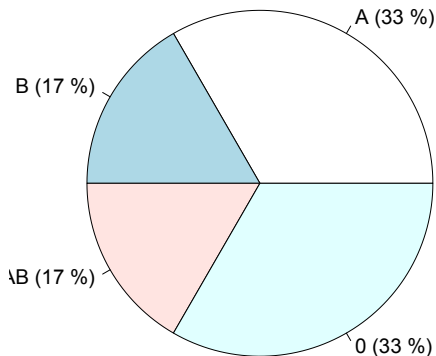
# Intermezzo

- určíme absolutní a relativní četnost krevní skupiny A
- $n_A = 4$ ;  $\pi_A = \frac{4}{12} = \frac{1}{3}$



# Koláčový diagram

- vhodný pro kvalitativní znaky k vyjádření četností jejich kategorií





# Klasická definice pravděpodobnosti

- je intuitivní a bude nám stačit
- pravděpodobnost jevu  $A$  je rovna podílu počtu případů  $m$ , které jsou jevu  $A$  příznivé, ku počtu  $n$  všech možným případů

$$P(A) = \frac{m}{n}$$

- nutným předpokladem je, že všechny případy mohou nastat stejně často

# Screeningové síto v medicíně

- způsob, jak plošně a včas diagnostikovat nemoci našeho zájmu
- měl by být levný
- síto složeno minimálně ze dvou po sobě jdoucích typů vyšetření
- každé z vyšetření má svou senzitivitu a specifitu

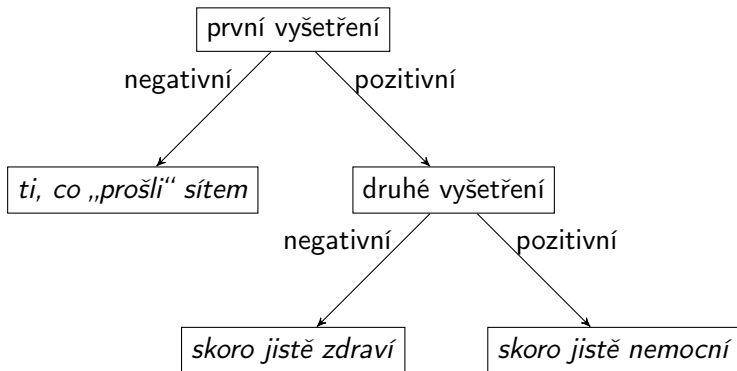
$$\text{senzitivita} = \frac{\# \text{ pozitivních}}{\# \text{ nemocných}}$$

$$\text{specifita} = \frac{\# \text{ negativních}}{\# \text{ zdravých}}$$

- první vyšetření by mělo být hodně senzitivní, druhé vyšetření by mělo být hodně specifické (viz další slide)

# Screeningové síto v medicíně

- prvním vyšetřením může být např. test na okultní krvácení, druhým vyšetření pak kolonoskopie





# Motivace

- ve výběru sto lidí je průměrná výška 175 cm a směrodatná odchylka je 10 cm
- jaká je s 95 % pravděpodobností průměrná výška populace?

---

# Pojem výběr

- vyšetřit celou populaci v praxi takřka nemožné
- nekonečně velké populace nelze celkově šetřit už z principu
- výběr := statistický soubor, obsahuje vybrané prvky z populace; je tedy podmnožinou populace
- výběr pořizujeme metodou náhodného, či záměrného výběru
- cílem získat reprezentativní výběr (vystihuje vlastnosti populace), nikoliv selektivní výběr

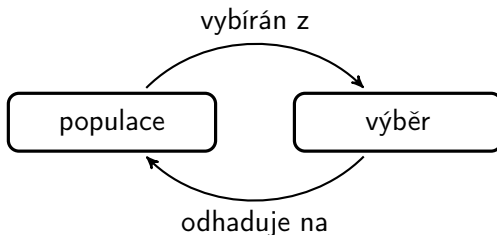
# Reprezentativní výběr

- takový výběr, z kterého je induktivními metodami možné usuzovat na vlastnosti „mateřské“ populace
- pořizujeme *záměrným*, či *náhodným* výběrem
  - *záměrný* výběr – opírá se o expertízu, zatížen subjektivitou
  - *náhodný* výběr – náhodné, nezávislé vybírání prvků populace do výběru



# Vztah populace a výběru

- z populace je vybírán výběr
- z charakteristik výběru jsou odhadovány charakteristiky populace (!)



# Bodový odhad statistického znaku

- předpokládáme, že charakteristická hodnota výběru (průměr, četnost) odpovídá populační hodnotě
- populační hodnota se pokládá rovna dané charakteristické hodnotě výběru
- např. „je-li četnost hypertoniků mezi dvaceti náhodnými pacienty 7, je i četnost hypertoniků v populaci  $\frac{7}{20} = 0,35 = 35 \%$ “
- s jakou „mírou jistoty“ jsme se „trefili“ do skutečné populační četnosti?
  - přirovnává se k lovu oštěpem

# Intervalový odhad statistického znaku

- (interval spolehlivosti, konfidenční interval)
- interval, ve kterém leží charakteristická hodnota populace s určitou pravděpodobností (spolehlivostí)
- např. např. „je-li četnost hypertoniků mezi dvaceti náhodnými pacienty 7, pak průměrná populační četnost hypertoniků leží s pravděpodobností 95 % intervalu (30; 40) %“
- s jakou „mírou jistoty“ jsme se „trefili“ do skutečné populační četnosti?
  - přirovnává se k lovu sítí

# Intervalový odhad statistického znaku

- máme-li výběr o rozsahu  $n$  s průměrem  $\bar{x}$  a směrodatnou odchylkou  $s$  daného znaku, pak populační průměr  $\mu$  daného znaku leží s 95 % pravděpodobností v intervalu

$$\mu \in \left( \bar{x} - 2\frac{s}{\sqrt{n}}; \bar{x} + 2\frac{s}{\sqrt{n}} \right)$$

- máme-li výběr o rozsahu  $n$  s relativní četností  $p$  daného znaku, pak populační relativní četnost  $\pi$  daného znaku leží s 95 % pravděpodobností v intervalu

$$\pi \in \left( \bar{p} - 2\sqrt{\frac{p(1-p)}{n}}; \bar{x} + 2\sqrt{\frac{p(1-p)}{n}} \right)$$

# Intermezzo

- ve výběru sto lidí je průměrná výška 175 cm a směrodatná odchylka je 10 cm
- jaká je s 95 % pravděpodobností průměrná výška populace?

# Intermezzo

- ve výběru sto lidí je průměrná výška 175 cm a směrodatná odchylka je 10 cm
- jaká je s 95 % pravděpodobností průměrná výška populace?
- 

$$\mu \in \left( \bar{x} - 2\frac{s}{\sqrt{n}}; \bar{x} + 2\frac{s}{\sqrt{n}} \right)$$

$$\mu \in \left( 175 - 2\frac{10}{\sqrt{100}}; 175 + 2\frac{10}{\sqrt{100}} \right)$$

$$\mu \in \left( 175 - 2\frac{10}{10}; 175 + 2\frac{10}{10} \right)$$

$$\mu \in (175 - 2; 175 + 2)$$

$$\mu \in (173; 177) \text{ [cm]}$$

# Princip testování hypotéz

- je založen na definování tzv. nulové hypotézy, kterou lze eventuálně vyvrátit nalezením významného protipříkladu
- nulovou hypotézou může být např. tvrzení, že průměrná výška v populaci je 171 cm
- protipříkladem je ve statistice myšlen dostatečně velký soubor hodnot, které jsou dostatečně „v rozporu“ s nulovou hypotézou
- protipříkladem může být např. výběr sto lidí, kde je průměrná výška 175 cm a směrodatná odchylka 10 cm

## Hladina významnosti

- předpokládejme, že nulová hypotéza platí; pak pravděpodobnost toho, že ji za její platnosti (chybně) zamítnu, by měla být co nejmenší a je nazývána  $p$ -hodnota či hladina signifikance
- hladina významnosti je tedy pravděpodobnost chyby (1. typu), proto by měla být co nejmenší
- je-li obvykle hladina významnosti  $\equiv p$ -hodnota  $\leq 0,05$ , lze již nulovou hypotézu zamítnout (riziko chyby prvního typu je malé)



# Intermezzo

- ve výběru sto lidí je průměrná výška 175 cm a směrodatná odchylka je 10 cm
- někdo tvrdí, že průměrná výška populace je 171 cm
- lze takovou hypotézu rozumně zamítnout?

# Intermezzo

- ve výběru sto lidí je průměrná výška 175 cm a směrodatná odchylka je 10 cm
- někdo tvrdí, že průměrná výška populace je 171 cm
- lze takovou hypotézu rozumně zamítnout?
- nulová hypotéza  $H_0 : \mu = 171$  [cm]

# Intermezzo

- ve výběru sto lidí je průměrná výška 175 cm a směrodatná odchylka je 10 cm
- někdo tvrdí, že průměrná výška populace je 171 cm
- lze takovou hypotézu rozumně zamítnout?
- nulová hypotéza  $H_0 : \mu = 171$  [cm]
- my ale díky předchozímu příkladu víme, že s 95 % pravděpodobností je  $\mu \in (173; 177)$  [cm]

# Intermezzo

- ve výběru sto lidí je průměrná výška 175 cm a směrodatná odchylka je 10 cm
- někdo tvrdí, že průměrná výška populace je 171 cm
- lze takovou hypotézu rozumně zamítnout?
- nulová hypotéza  $H_0 : \mu = 171$  [cm]
- my ale díky předchozímu příkladu víme, že s 95 % pravděpodobností je  $\mu \in (173; 177)$  [cm]
- pravděpodobnost chyby (1. typu) při zamítnutí nulové hypotézy je tak menší než  $100 \% - 95 \% = 5 \%$

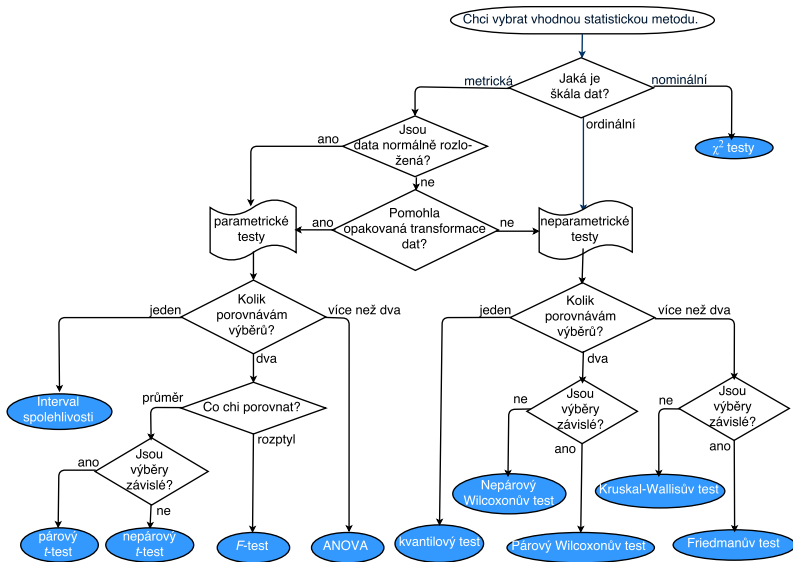
# Intermezzo

- ve výběru sto lidí je průměrná výška 175 cm a směrodatná odchylka je 10 cm
- někdo tvrdí, že průměrná výška populace je 171 cm
- lze takovou hypotézu rozumně zamítnout?
- nulová hypotéza  $H_0 : \mu = 171$  [cm]
- my ale díky předchozímu příkladu víme, že s 95 % pravděpodobností je  $\mu \in (173; 177)$  [cm]
- pravděpodobnost chyby (1. typu) při zamítnutí nulové hypotézy je tak menší než  $100 \% - 95 \% = 5 \%$
- nulovou hypotézu tak lze zamítnout

# Testy hypotéz

- předchozí úvahy ale obvykle není nutné pokaždé provádět, existují zavedené algoritmy, tzv. testy hypotéz, které vrací pouze hladinu signifikance ( $p$ -hodnotu)
- je-li  $p$ -hodnota  $\leq 0,05$ , zamítáme nulovou hypotézu

# Testy hypotéz



# Literatura



Karel Zvára. *Biostatistika*. Praha: Karolinum, 2003. ISBN: 978-80-246-0739-9.



Jana Zvárová. *Základy statistiky pro biomedicínské obory*. Praha: Karolinum, 2016. ISBN: 978-80-246-3416-6.



Děkuji za pozornost!

lubomir.stepanek@lf1.cuni.cz

lubomir.stepanek@fbmi.cvut.cz