

# Základy statistiky

## Explorační analýza dat

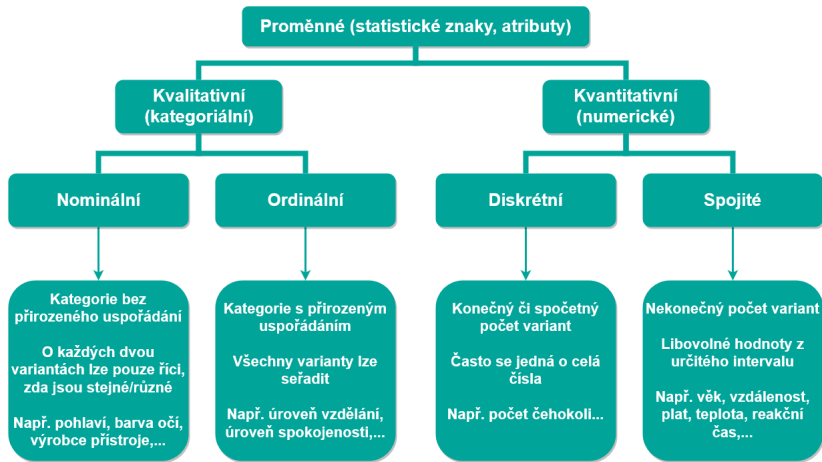
Adéla Vrtková, Martina Litschmannová

[adela.vrtkova@vsb.cz](mailto:adela.vrtkova@vsb.cz), [martina.litschmannova@vsb.cz](mailto:martina.litschmannova@vsb.cz)

 **VŠB** TECHNICKÁ  
UNIVERZITA  
OSTRAVA

FAKULTA  
ELEKTROTECHNIKY  
A INFORMATIKY

KATEDRA  
APLIKOVANÉ  
MATEMATIKY





- tabulky četností
  - absolutní a relativní četnosti
  - popř. kumulativní abs. a rel. četnosti pro ordinální proměnnou
- sloupcové grafy
- výšečové grafy

## Tabulka rozdělení četnosti

Varianty	Absolutní četnosti	Relativní četnosti
$x_i$	$n_i$	$p_i$
$x_1$	$n_1$	$p_1 = n_1/n$
$x_2$	$n_2$	$p_2 = n_2/n$
$\vdots$	$\vdots$	$\vdots$
$x_k$	$n_k$	$p_k = n_k/n$
<b>Celkem</b>	$n = \sum_i n_i$	<b>1</b>

**Tabulka rozdělení četnosti**

<b>Věková kategorie</b>	<b>Absolutní četnosti</b>	<b>Relativní četnosti (%)</b>
pod 35 let	77	37,37864
35 až 50 let	85	41,26214
nad 50 let	44	21,35922
<b>Celkem</b>	<b>206</b>	<b>100</b>

**Tabulka rozdělení četnosti**

<b>Věková kategorie</b>	<b>Absolutní četnosti</b>	<b>Relativní četnosti (%)</b>
pod 35 let	77	37,37864
35 až 50 let	85	41,26214
nad 50 let	44	21,35922
<b>Celkem</b>	<b>206</b>	<b>100</b>



Jak zaokrouhlit relativní četnost?

**Tabulka rozdělení četnosti**

<b>Věková kategorie</b>	<b>Absolutní četnosti</b>	<b>Relativní četnosti (%)</b>
pod 35 let	77	37,4
35 až 50 let	85	41,3
nad 50 let	44	21,4
<b>Celkem</b>	<b>206</b>	<b>100,1</b>

**Tabulka rozdělení četnosti**

<b>Věková kategorie</b>	<b>Absolutní četnosti</b>	<b>Relativní četnosti (%)</b>
pod 35 let	77	37,4
35 až 50 let	85	41,3
nad 50 let	44	21,4
<b>Celkem</b>	<b>206</b>	<b>100,1</b>



Pozor na zaokrouhlovací chybu!

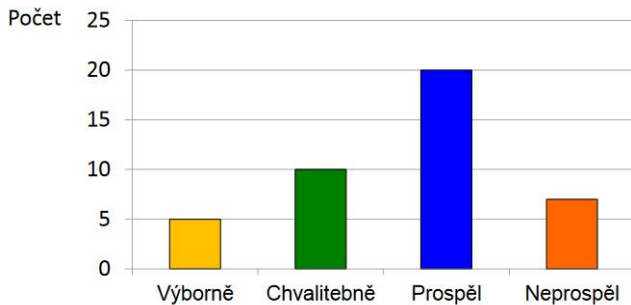


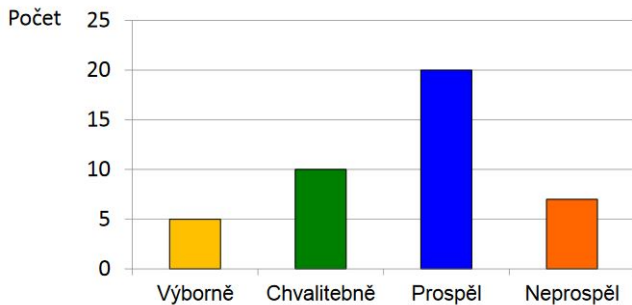
**Tabulka rozdělení četnosti**

<b>Věková kategorie</b>	<b>Absolutní četnosti</b>	<b>Relativní četnosti (%)</b>
pod 35 let	77	37,4
35 až 50 let	85	41,3
nad 50 let	44	21,3
<b>Celkem</b>	<b>206</b>	<b>100,0</b>

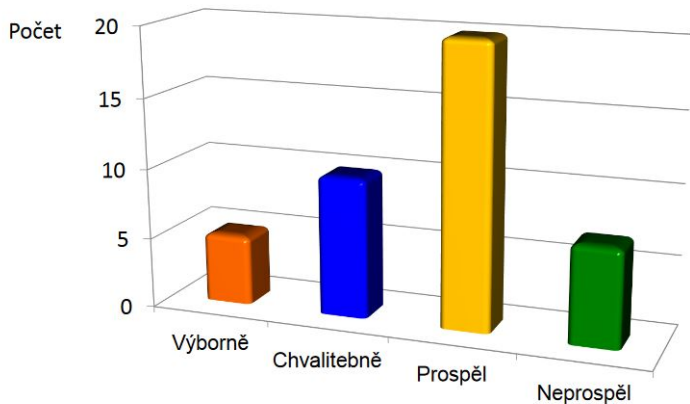


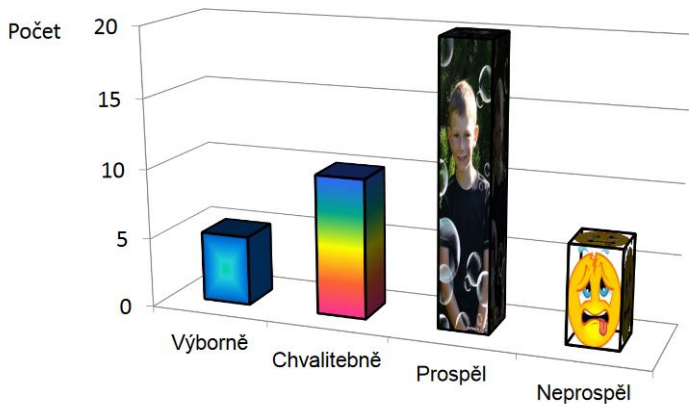
- sloupcové grafy
- výsečové grafy

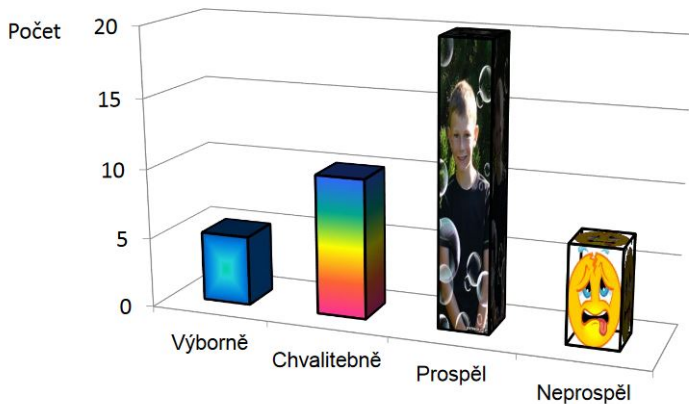




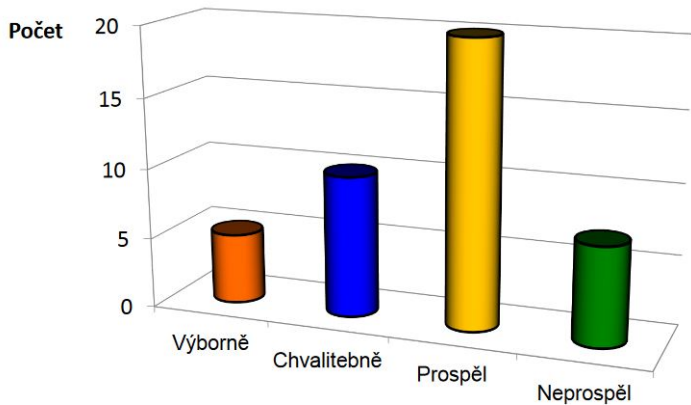
Návrhy na vylepšení?



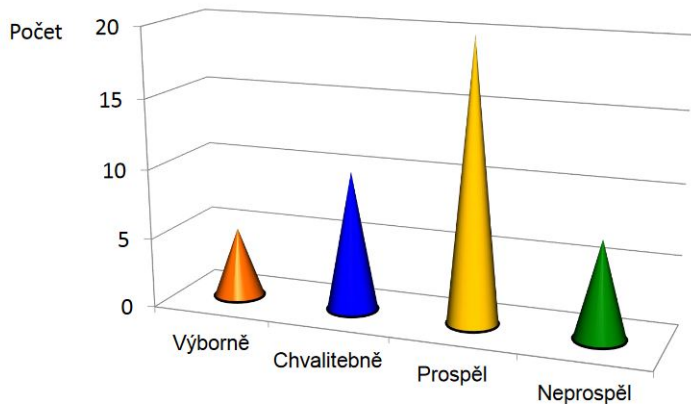


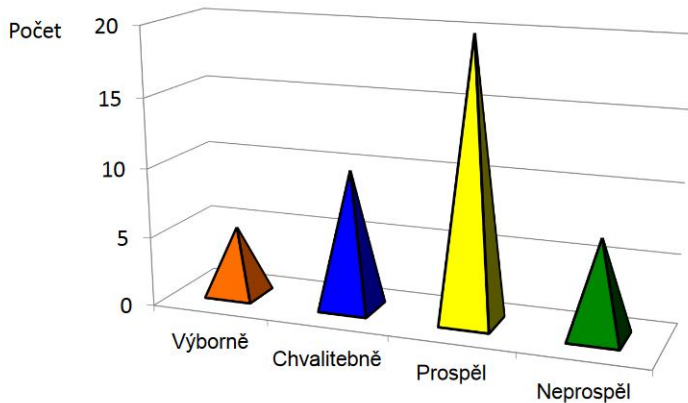


Nezapomínejme, že někdy „méně je více“.









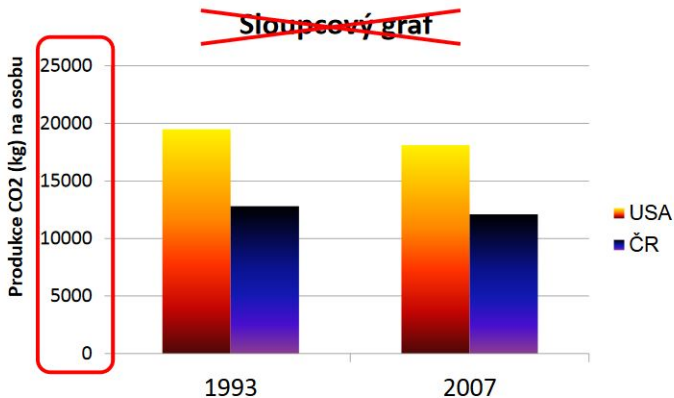


- vyhnout se 3D grafům

„3D is great when it comes to James Cameron movies, but your data does not need to jump out at your viewers“  
(zdroj: visage.co)



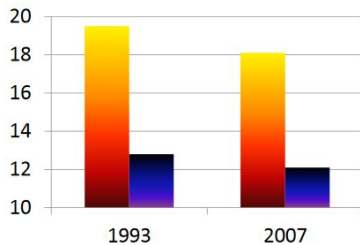
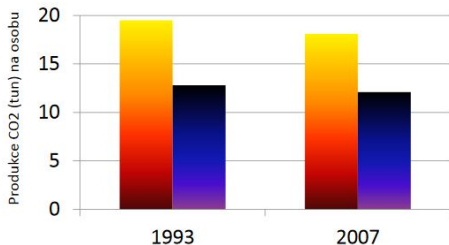
- neefektivní nuly a nadbytečné názvy a popisy



(zdroj dat: [http://en.wikipedia.org/wiki/List\\_of\\_countries\\_by\\_carbon\\_dioxide\\_emissions\\_per\\_capita](http://en.wikipedia.org/wiki/List_of_countries_by_carbon_dioxide_emissions_per_capita))



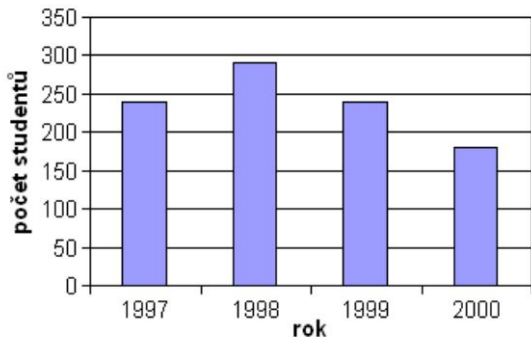
■ pozor na efekt změny rozsahu osy y





Určete pravdivost tvrzení:

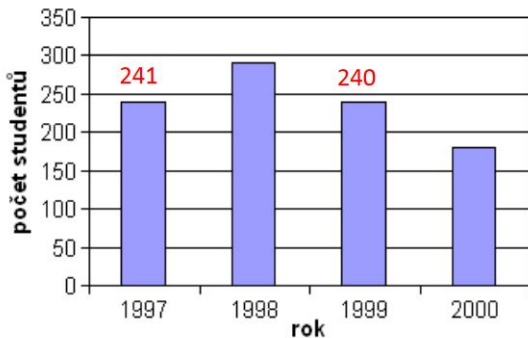
V žádných dvou letech nebyl počet studentů stejný.

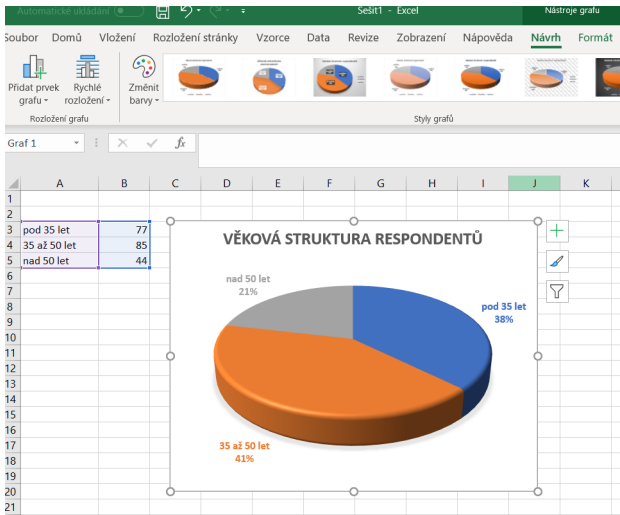




Určete pravdivost tvrzení:


V žádných dvou letech nebyl počet studentů stejný.

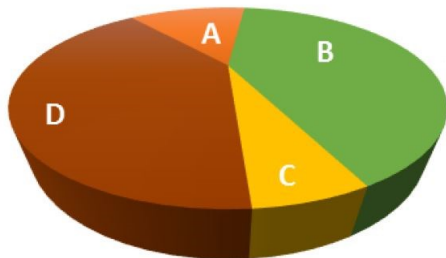







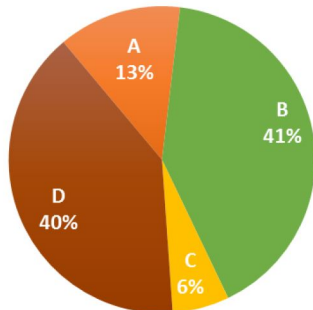
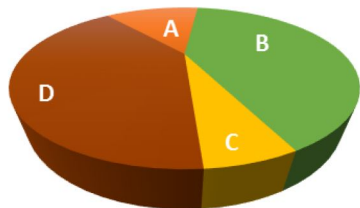


 Jaký je poměr mezi velikostí výsečí A a C? Jaký je poměr mezi velikostí výsečí B a D?





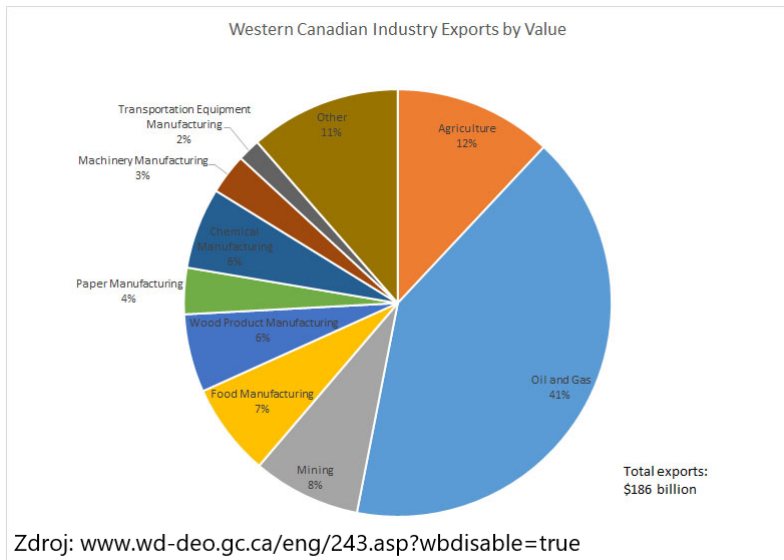
 Jaký je poměr mezi velikostí výsečí A a C? Jaký je poměr mezi velikostí výsečí B a D?





- citace z nápovědy ke statistickému softwaru R:

"Pie charts are a very bad way of displaying information. The eye is good at judging linear measures and bad at judging relative areas. A bar chart or dot chart is a preferable way of displaying this type of data."





- číselné charakteristiky
- krabicové grafy
- histogramy



- míry polohy
- míry variability
- míry tvaru



- míry polohy
- míry variability
- míry tvaru



Všechny míry lze vypočítat pro výběr i pro populaci (pokud jsou k dispozici data).



- Míry polohy, také označované jako míry centrální tendence, cílí nejčastěji na odhad jakéhosi „středu“.
- Nejčastěji používané míry polohy:
  - aritmetický průměr,
  - medián,
  - modus,
  - minimum a maximum,
  - kvantily (kvartily, decily, percentily).





$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$



- V malé vesnici někde v Americe žije 6 lidí, jejichž roční plat (v USD) je uveden níže.

25 000, 27 000, 29 000, 35 000, 37 000, 38 000.

Určete průměrný plat obyvatel této vesnice.



- V malé vesnici někde v Americe žije 6 lidí, jejichž roční plat (v USD) je uveden níže.

25 000, 27 000, 29 000, 35 000, 37 000, 38 000.

Určete průměrný plat obyvatel této vesnice.

**31 830**



- Do vesnice se přistěhoval Bill Gates, jehož roční příjem je 40 milionů dolarů.

25 000, 27 000, 29 000, 35 000, 37 000, 38 000,  
40 000 000.

Jaký je průměrný plat obyvatel této vesnice nyní?



- Do vesnice se přistěhoval Bill Gates, jehož roční příjem je 40 milionů dolarů.

25 000, 27 000, 29 000, 35 000, 37 000, 38 000,  
40 000 000.

Jaký je průměrný plat obyvatel této vesnice nyní?

**5 741 571**



- Do vesnice se přistěhoval Bill Gates, jehož roční příjem je 40 milionů dolarů.

25 000, 27 000, 29 000, 35 000, 37 000, 38 000,  
40 000 000.

Jaký je průměrný plat obyvatel této vesnice nyní?

**5 741 571**



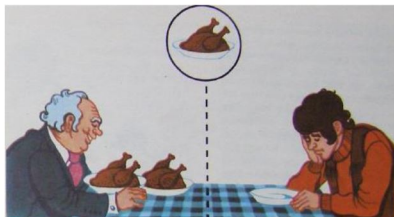
Aritmetický průměr je velice citlivý na odlehlé hodnoty!



Matematik, fyzik a statistik střílejí na terč. Matematik střílí první a mine cíl kousek doleva. Fyzik střílí po něm a také mine cíl tentokrát kousek doprava. Statistik se začne radovat: „Hurá! Trefili jsme cíl!“.



Matematik, fyzik a statistik střílejí na terč. Matematik střílí první a mine cíl kousek doleva. Fyzik střílí po něm a také mine cíl tentokrát kousek doprava. Statistik se začne radovat: „Hurá! Trefili jsme cíl!“.



*Zdroj: SWOBODA, Helmut. Moderní statistika., 1977.*





$100p\%$  kvantil  $\tilde{x}_p$

- $100p\%$  hodnot (dat) je menší nebo rovno hodnotě  $\tilde{x}_p$
- např. 20 % hodnot je menší nebo rovno **20% kvantilu**, tudíž 80 % hodnot je nad ním



**Medián** je 50% kvantil.



## ■ Kvantily

- dolní kvartil = 25% kvantil ( $\tilde{x}_{0,25}$ )
- medián = 50% kvantil ( $\tilde{x}_{0,50}$ )
- horní kvartil = 75% kvantil ( $\tilde{x}_{0,75}$ )

## ■ Decily

- 1. decil = 10% kvantil ( $\tilde{x}_{0,10}$ )
- 2. decil = 20% kvantil ( $\tilde{x}_{0,20}$ )
- $\vdots$
- 9. decil = 90% kvantil ( $\tilde{x}_{0,90}$ )

## ■ Percentily

- 1. percentil = 1% kvantil ( $\tilde{x}_{0,01}$ )
- 2. percentil = 2% kvantil ( $\tilde{x}_{0,02}$ )
- $\vdots$
- 99. percentil = 99% kvantil ( $\tilde{x}_{0,99}$ )

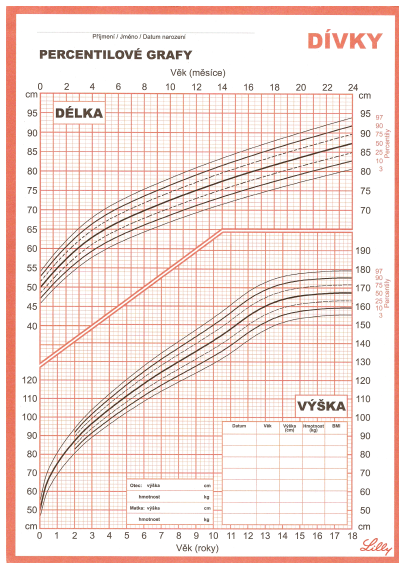


Jaké hodnoty reprezentují 0% a 100% kvantil?



Jaké hodnoty reprezentují 0% a 100% kvantil?

**minimum a maximum**



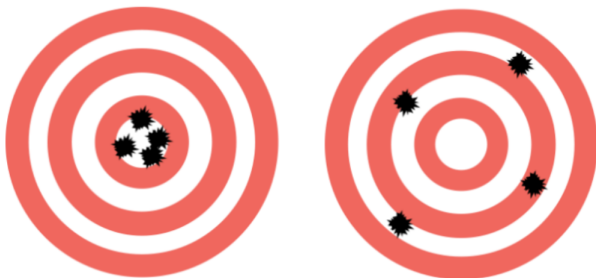


„80 % nemovitostí prodáme do 2 měsíců“

Zdroj: výloha Fincentrum Reality, 28. října 813/248, 70900 Ostrava



- Charakteristiky hodnotící rozptýlenost hodnot statistického souboru kolem nějaké míry polohy.
- Nejčastěji používané míry variability:
  - rozptyl
  - směrodatná odchylka
  - variační koeficient
  - rozpětí
  - mezikvartilové rozpětí

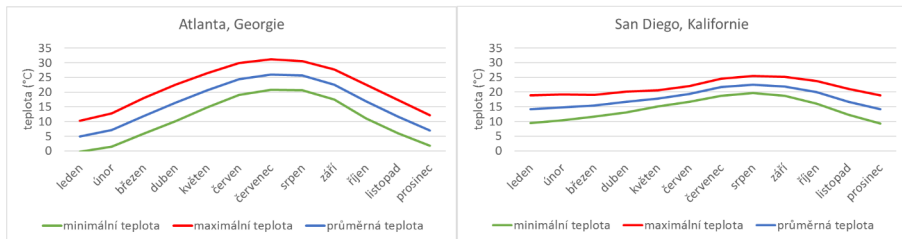


Zdroj: blackswanfarming.com





- Průměrná teplota v Atlantě je  $15,5^{\circ}\text{C}$ .
- Průměrná teplota v San Diegu je  $16,5^{\circ}\text{C}$ .





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



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



Jednotky rozptylu jsou druhou mocninou jednotek analyzované proměnné!



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



- Čebyševova nerovnost:  $\forall k > 0 : P(\mu - k\sigma < X < \mu + k\sigma) > 1 - \frac{1}{k^2}$

$k$	$P(\mu - k\sigma < X < \mu + k\sigma)$
1	$> 0$
2	$> 0,75$
3	$> 0,89$

- $\mu \cong \bar{x}$
- $\sigma \cong s$



- Používá se většinou pro proměnné nabývající nezáporných hodnot.
- Umožňuje srovnání variability proměnných, které mají různé jednotky.

$$V = \frac{s}{|\bar{x}|}$$

- Obvykle vyjádřen v procentech. Čím je nižší, tím homogennější soubor.

$$V = \frac{s}{|\bar{x}|} \cdot 100 \quad (\%)$$

- Opatrně, je-li průměr blízký nule.



$$\text{Rozpětí} = \text{Max} - \text{Min}$$



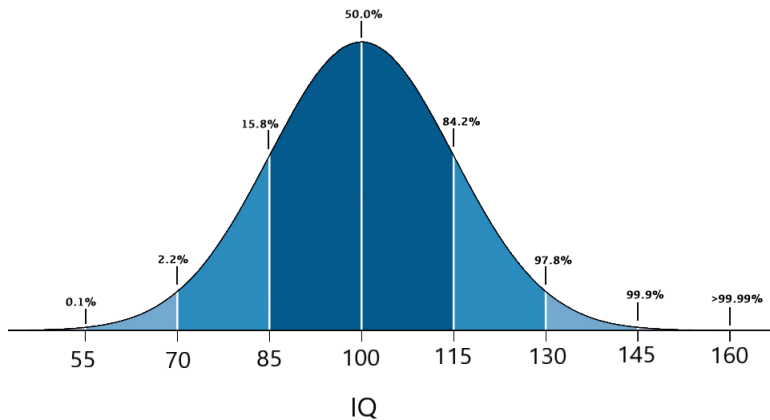
- Využívá se např. při identifikaci odlehlých pozorování.

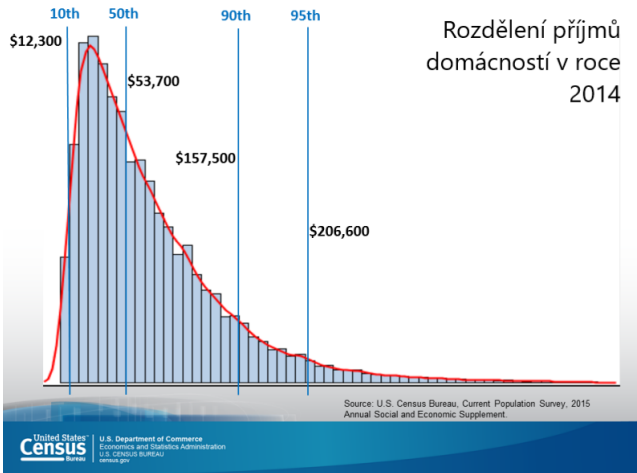
$$IQR = \tilde{x}_{0,75} - \tilde{x}_{0,25}$$

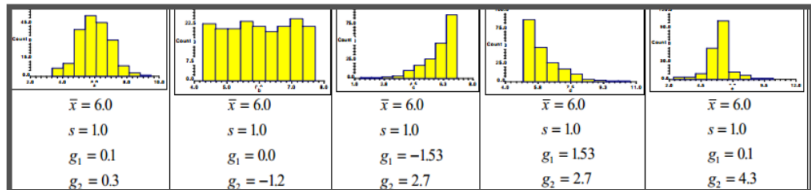




- Míry tvaru popisují symetrii dat a koncentraci dat kolem průměru.
- Slouží ke srovnání rozdělení dat s rozdělením normálním.
- Nejčastěji používané míry tvaru:
  - šikmost
  - špičatost





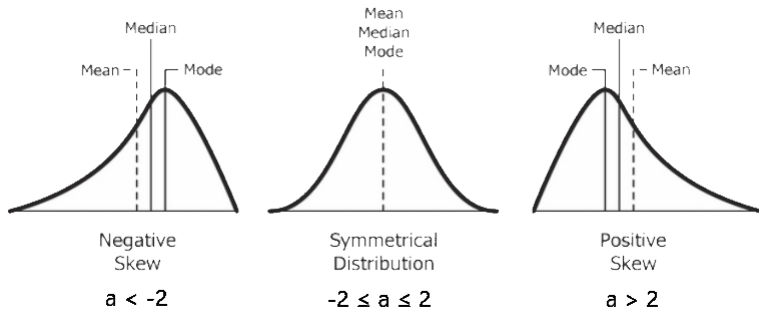


Zdroj: TVRDÍK, J.: Základy matematické statistiky, Ostravská univerzita, 2008

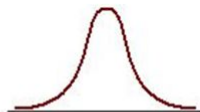
- Všechny ukázky mají stejný průměr a směrodatnou odchylku. Na první pohled je však vidět různý tvar rozdělení. Právě tyto rozdíly nám číselně dokážou popsat míry tvaru - šikmost a špičatost (na obrázku značeno  $g_1$  a  $g_2$ ).



$$a = \frac{n}{(n-1)(n-2)} \cdot \frac{\sum_{i=1}^n (x_i - \bar{x})^3}{s^3}$$



$$b = \frac{n(n+1)}{(n-1)(n-2)(n-3)} \cdot \frac{\sum_{i=1}^n (x_i - \bar{x})^4}{s^4} - 3 \frac{(n-1)^2}{(n-2)(n-3)}$$

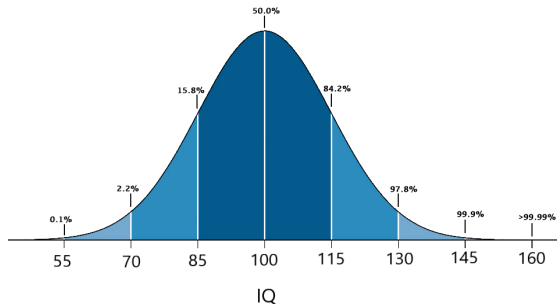
 $b < -2$  $-2 \leq b \leq 2$  $b > 2$



Míry tvaru mohou být v různých softwarech různě definovány. Je důležité si definici vždy v nápovědě zkontrolovat, abychom výsledky správně interpretovali.



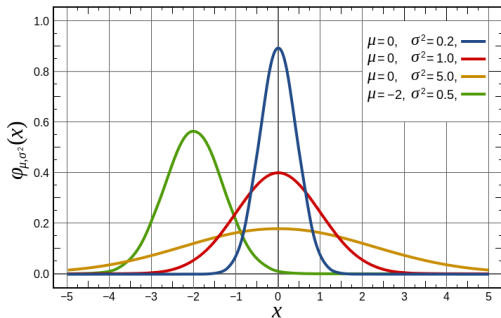
- nazýváno také Gaussovské
- jedná se o jakési „chování“ dat, které má výhodné vlastnosti



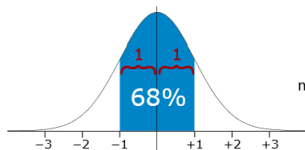




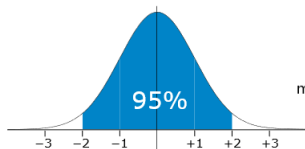
- $N(\mu, \sigma^2)$  -  $\mu$  (průměr),  $\sigma^2$  (rozptyl)
- výběr pocházející z normálního rozdělení by měl mít šikmost a špičatost přibližně nulovou
- průměr = modus = medián



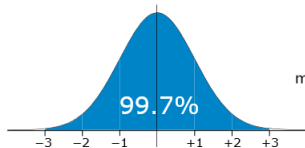
[www.statisticshowto.com/probability-and-statistics/normal-distributions/](http://www.statisticshowto.com/probability-and-statistics/normal-distributions/)



**asi 68 % hodnot** se nachází ve vzdálenosti menší než **1 směrodatná odchylka** od průměru



**asi 95 % hodnot** se nachází ve vzdálenosti menší než **2 směrodatné odchylky** od průměru



**asi 99,7 % hodnot** se nachází ve vzdálenosti menší než **3 směrodatné odchylky** od průměru

[www.mathsisfun.com/data/standard-normal-distribution.html](http://www.mathsisfun.com/data/standard-normal-distribution.html)



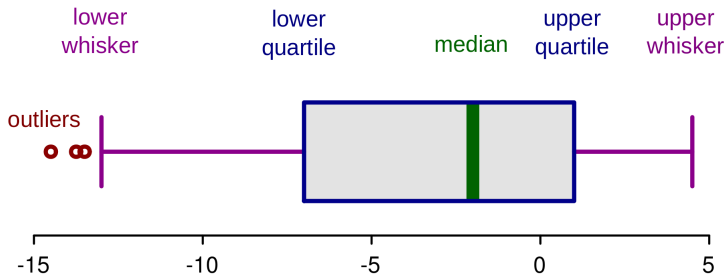
- Hodnoty, které se **mimořádně liší od ostatních hodnot** analyzované proměnné a tím ovlivňují např. vypovídací hodnotu průměru.
- Můžou se vyskytnout zcela **náhodně, chybou měření, lidskou chybou** apod.
- Známe-li příčinu odlehlosti a předpokládáme-li, že již nenastane, jsme oprávněni tato pozorování vyloučit.
- V ostatních případech je nutno zvážit, zda se vyloučením odlehlých pozorování nepřipravíme o důležité informace.
  - Metoda vnitřních/vnějších hradeb
  - Chauvenetovo kritérium
  - Mahalanobisova vzdálenost

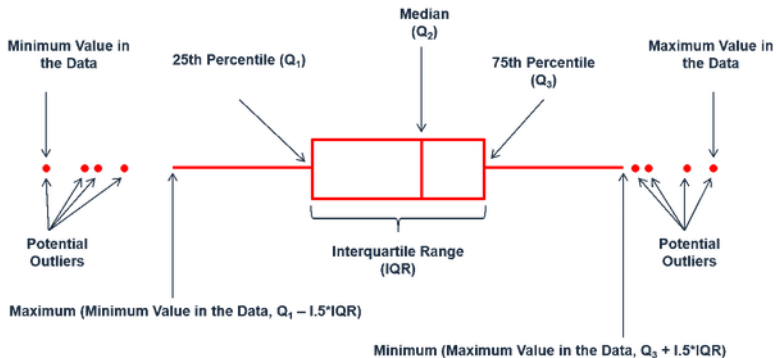


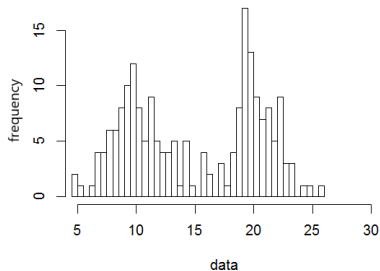
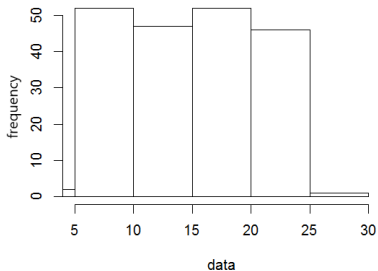
Je-li pozorování  
**menší než  $\tilde{x}_{0,25} - 1,5 \cdot \text{IQR}$**   
nebo  
**větší než  $\tilde{x}_{0,75} + 1,5 \cdot \text{IQR}$ ,**  
mělo by být označeno za odlehlé pozorování.



- krabicové grafy (= boxploty)
- histogramy
- QQ-grafy



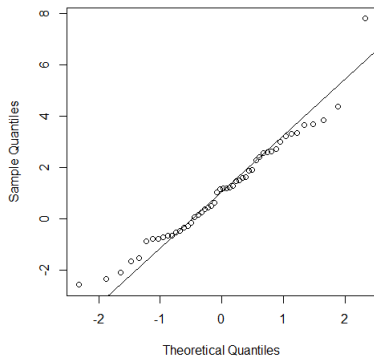




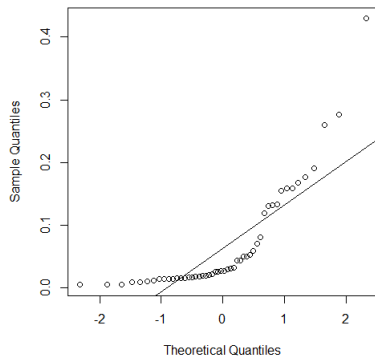




Normal Q-Q Plot



Normal Q-Q Plot





Jak zaokrouhlit míry polohy? Jak zaokrouhlit směrodatnou odchylku?

BMI			
rozsah	průměr	medián	sm. odchylka
30	25,337	24,165	5,147



Zaokrouhlíme nahoru směrodatnou odchylku na určitý počet platných cifer podle pravidla (viz níže).

rozsah výběru	$\langle 2; 10 \rangle$	$(10; 30)$	$(30; 2000)$	$(2000; \dots)$
počet platných cifer	1	2	3	4

BMI			
rozsah	průměr	medián	sm. odchylka
30	25,337	24,165	5,147



Zaokrouhlíme nahoru směrodatnou odchylku na určitý počet platných cifer podle pravidla (viz níže).

rozsah výběru	$\langle 2; 10 \rangle$	$(10; 30)$	$(30; 2000)$	$(2000; \dots)$
počet platných cifer	1	2	3	4

BMI			
rozsah	průměr	medián	sm. odchylka
30	25,337	24,165	<b>5,2</b>



Zaokrouhlíme nahoru směrodatnou odchylku na určitý počet platných cifer podle pravidla a **míry polohy jsou zaokrouhleny na stejný počet desetinných míst, resp. na stejný řád.**

BMI			
rozsah	průměr	medián	sm. odchylka
30	25,337	24,165	<b>5,2</b>



Zaokrouhlíme nahoru směrodatnou odchylku na určitý počet platných cifer podle pravidla a **míry polohy jsou zaokrouhleny na stejný počet desetinných míst, resp. na stejný řád.**

BMI			
rozsah	průměr	medián	sm. odchylka
30	<b>25,3</b>	<b>24,2</b>	<b>5,2</b>



Zaokrouhlíme nahoru směrodatnou odchylku na určitý počet platných cifer podle pravidla a **míry polohy jsou zaokrouhleny na stejný počet desetinných míst, resp. na stejný řád.**

BMI			
rozsah	průměr	medián	sm. odchylka
30	<b>25,3</b>	<b>24,2</b>	<b>5,2</b>



Toto pravidlo lze aplikovat na všechny míry polohy s výjimkou minima a maxima, kde preferujeme hodnoty přímo z datového souboru.