

Základy statistiky

Statistika kolem nás

Základní pojmy statistické analýzy dat

Adéla Vrtková, Martina Litschmannová

adela.vrtkova@vsb.cz, martina.litschmannova@vsb.cz





- odvětví matematiky zabývající se sběrem, analýzou, interpretací a prezentací dat (Merriam-Webster.com)



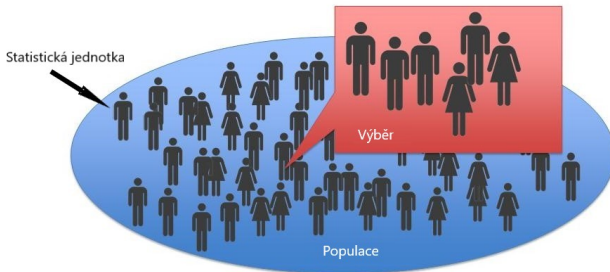
- odvětví matematiky zabývající se sběrem, analýzou, interpretací a prezentací dat (Merriam-Webster.com)
- věda, která se snaží zkoumat reálná data a s pomocí teorie pravděpodobnosti se tato data snaží popisovat (Matematika.cz)

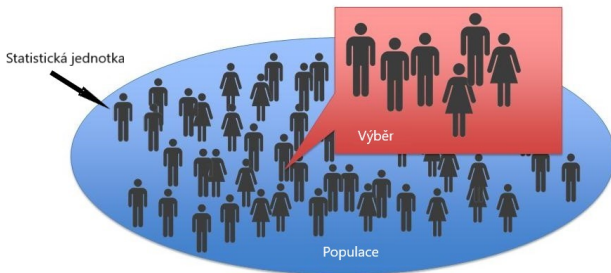


- odvětví matematiky zabývající se sběrem, analýzou, interpretací a prezentací dat (Merriam-Webster.com)
- věda, která se snaží zkoumat reálná data a s pomocí teorie pravděpodobnosti se tato data snaží popisovat (Matematika.cz)
- přesný součet nepřesných čísel (Poznámka k lidové definici statistiky, časopis Vesmír 5/1998)



- odvětví matematiky zabývající se **sběrem, analýzou, interpretací a prezentací dat** (Merriam-Webster.com)
- věda, která se snaží zkoumat reálná data a s pomocí teorie pravděpodobnosti se tato data snaží popisovat (Matematika.cz)
- přesný součet nepřesných čísel (Poznámka k lidové definici statistiky, časopis Vesmír 5/1998)





- **populace** - všechny statistické jednotky
- **výběr** - podmnožina populace



Představte si, že chcete provést průzkum spokojenosti studentů VŠB-TUO se službami knihovny. Asi není možné se zeptat úplně všech studentů univerzity, resp. těch, kteří služby knihovny využívají. Jak tedy provést výběr studentů do průzkumu?



Představte si, že chcete provést průzkum spokojenosti studentů VŠB-TUO se službami knihovny. Asi není možné se zeptat úplně všech studentů univerzity, resp. těch, kteří služby knihovny využívají. Jak tedy provést výběr studentů do průzkumu?

- Způsoby provedení výběru se zabývá **teorie výběrového šetření** (Steven K. Thompson: Sampling).



Představte si, že chcete provést průzkum spokojenosti studentů VŠB-TUO se službami knihovny. Asi není možné se zeptat úplně všech studentů univerzity, resp. těch, kteří služby knihovny využívají. Jak tedy provést výběr studentů do průzkumu?

- Způsoby provedení výběru se zabývá **teorie výběrového šetření** (Steven K. Thompson: Sampling).
- Ta úzce souvisí s **metodami sběru dat**.



Představte si, že chcete provést průzkum spokojenosti studentů VŠB-TUO se službami knihovny. Asi není možné se zeptat úplně všech studentů univerzity, resp. těch, kteří služby knihovny využívají. Jak tedy provést výběr studentů do průzkumu?

- Způsoby provedení výběru se zabývá **teorie výběrového šetření** (Steven K. Thompson: Sampling).
- Ta úzce souvisí s **metodami sběru dat**.
- Cílem je získání **reprezentativního výběru**, což je takový výběr, který má stejné vlastnosti (strukturu) jako populace.



Prostý náhodný výběr (bez vracení / s vracením)

- los z osudí, generátor pseudonáhodných čísel

Systematický výběr

- výběr z populace podle pevně zvoleného kroku - každou třetí jednotku, každou desátou, apod.

Stratifikovaný výběr

- rozklad populace na disjunktní oblasti (strata) dle určitého kritéria; v každé oblasti pak prostý náhodný výběr

Skupinkový výběr

- populace rozdělena do homogenních skupin, prostý náhodný nebo systematický výběr několika skupin, sledování celých vybraných skupin

Vícestupňový výběr

- analogie skupinkového výběru, u vybraných skupin podskupiny atd.

Výběr s nesterjnými pravděpodobnostmi

- statistické jednotky mají různou pravděpodobnost zahrnutí do výběru



Kvótní výběr

- snaha naplnit předem dané kvóty (věkové kategorie, pohlaví apod.), nutná znalost struktury populace

Snowball technika

- výběr malého souboru jednotek splňující kritéria, skrze ně získání kontaktů na další

Účelový výběr

- založen zcela na úsudku výzkumníka, nelze očekávat zobecnitelné výsledky

Příležitostný výběr

- „všichni, kteří jsou po ruce“

Samovýběr

- účast na základě rozhodnutí statistické jednotky



- rozhovory
- dotazníková šetření a ankety
- pozorování
- ohniskové skupiny (focus groups)
- případové studie
- sekundární data, analýza dokumentů
- experimenty, měření



- rozhovory
- **dotazníková šetření a ankety**
- pozorování
- ohniskové skupiny (focus groups)
- případové studie
- sekundární data, analýza dokumentů
- **experimenty, měření**

prostý náhodný výběr

systematický výběr

stratifikovaný výběr

skupinkový výběr

vícetupňový výběr

výběr s nestejnými p-stmi

kvótní výběr

snowball technika

účelový výběr

příležitostný výběr

samovýběr

rozhovory

dotazníková šetření

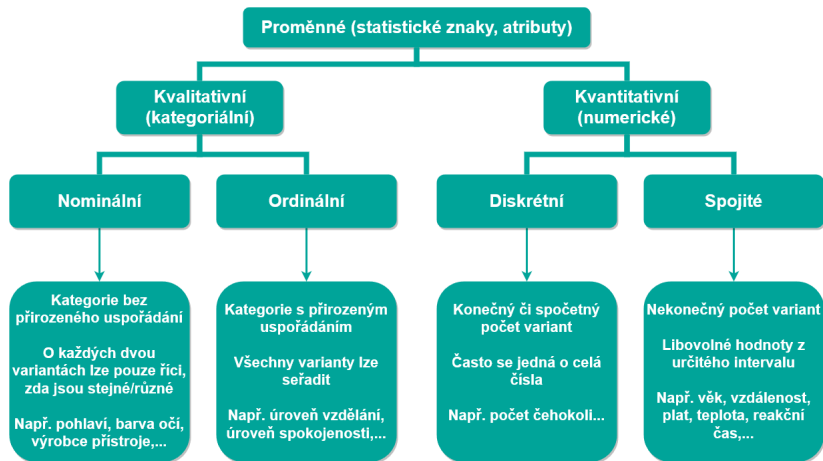
pozorování

ohniskové skupiny

případové studie

sekundární data

experimenty, měření





Provedli jsme výběr statistických jednotek, sebrali jsme data, jak je ale zapsat?



Provedli jsme výběr statistických jednotek, sebrali jsme data, jak je ale zapsat?

- 1 řádek = 1 statistická jednotka (člověk, měření,...)
- 1 sloupec = 1 proměnná



Provedli jsme výběr statistických jednotek, sebrali jsme data, jak je ale zapsat?

- 1 řádek = 1 statistická jednotka (člověk, měření,...)
- 1 sloupec = 1 proměnná

	Proměnná 1	...	Proměnná p
1	x_{11}	...	x_{1p}
\vdots	\vdots	\vdots	\vdots
n	x_{n1}	...	x_{np}



- 1 řádek = 1 statistická jednotka (člověk, měření,...)
- 1 sloupec = 1 proměnná

ID	věk	hmotnost_před (kg)	hmotnost_po (kg)	výška (cm)	pohlaví
1	23	73,4	70,0	181	M
2	22	79,2	78,6	189	M
3	22	71,9	69,9	167	Ž
4	26	66,3	63,8	165	Ž
5	24	80,3	75,6	170	Ž
⋮	⋮	⋮	⋮	⋮	⋮



- jakákoli jiná struktura dat
- pro tyto jiné struktury budeme používat termín *tabulka* nebo *datová tabulka*

muži				ženy			
věk	hmotnost_před (kg)	hmotnost_po (kg)	výška (cm)	věk	hmotnost_před (kg)	hmotnost_po (kg)	výška (cm)
23	83,4	80,0	175	25	73,4	70,0	165
22	99,2	88,6	182	23	69,2	68,6	169
22	81,9	79,9	190	21	71,9	69,9	172
26	89,3	83,8	182	27	69,3	63,8	160
24	88,3	83,6	180	20	78,3	73,6	170
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮



Konečně máme sesbíraná a správně zapsaná data. Co s nimi?

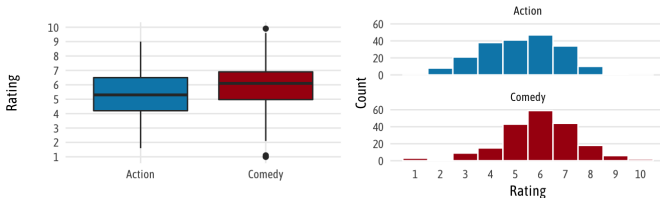
- explorační analýza dat
- aplikace metod statistické indukce



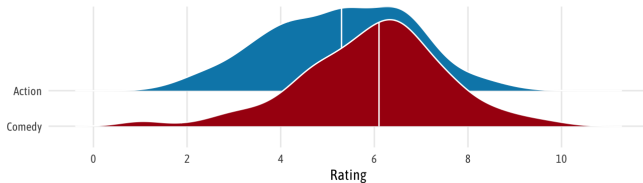
= popisná statistika, cílí na **prvotní průzkum dat**, nalezení **chyb, odlehlých pozorování**, získání přehledu o tom, s jakými daty pracujeme. Díky závěrům explorační analýzy můžeme také **zformulovat statistické hypotézy** nebo dopředu **odhadnout odpověď** na některé předem definované výzkumné otázky. Zároveň je nutnou součástí **analýzy předpokladů** statistických metod.

Do comedies get higher ratings than action movies?

Sample of 400 movies from IMDB



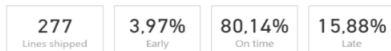
White line shows median rating



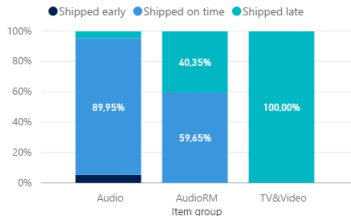
Zdroj: www.andrewheiss.com/blog/2019/01/29/diff-means-half-dozen-ways/



Shipping performance



Shipped by product



Shipped by site / warehouse



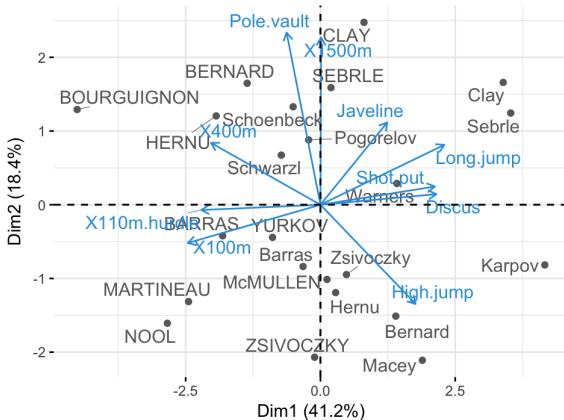
Zdroj: ax-dynamics.com/article/4-more-power-bi-content-packs-for-2017



Z těch nejznámějších metod spadajících pod tuto oblast uvedeme např. **testování hypotéz, regresní modely, metodu hlavních komponent, shlukovou analýzu nebo analýzu přežití**. Jedná se prakticky o všechny metody, kdy jsou **na základě analýzy výběru vyvozovány závěry o populaci**.



Metoda hlavních komponent



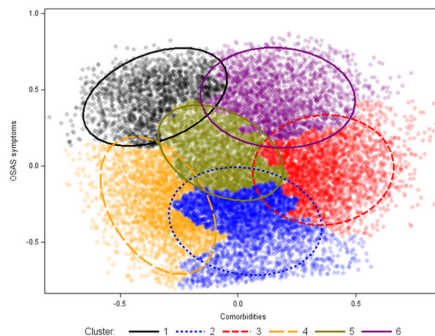
Zdroj: www.sthda.com/english/articles/31-principal-component-methods-in-r-practical-guide/112-pca-principal-component-analysis-essentials/

Shluková analýza

Obstructive Sleep Apnea: A Cluster Analysis at Time of Diagnosis

Representation of six clusters after ascending hierarchical clustering analysis.

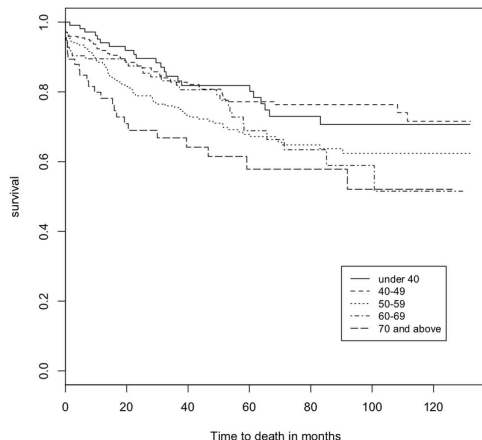
Axes correspond to individual coordinates for the two main dimensions of the multiple correspondence analysis. Cluster 1: the young symptomatic. Cluster 2: the old obese. Cluster 3: the multi-disease (MD) old obese. Cluster 4: the young snorers. Cluster 5: the drowsy obese. Cluster 6: the MD obese symptomatic.



Zdroj: BAILLY, Sébastien, et al. Obstructive sleep apnea: a cluster analysis at time of diagnosis. PLoS One, 2016, 11.6: e0157318.



Analýza přežití



Zdroj: KATANODA, Kota; MATSUDA, Tomohiro. Five-year relative survival rate of breast cancer in the USA, Europe and Japan. Japanese journal of clinical oncology, 2014, 44.6: 611-611.