

1. Statistická analýza dat

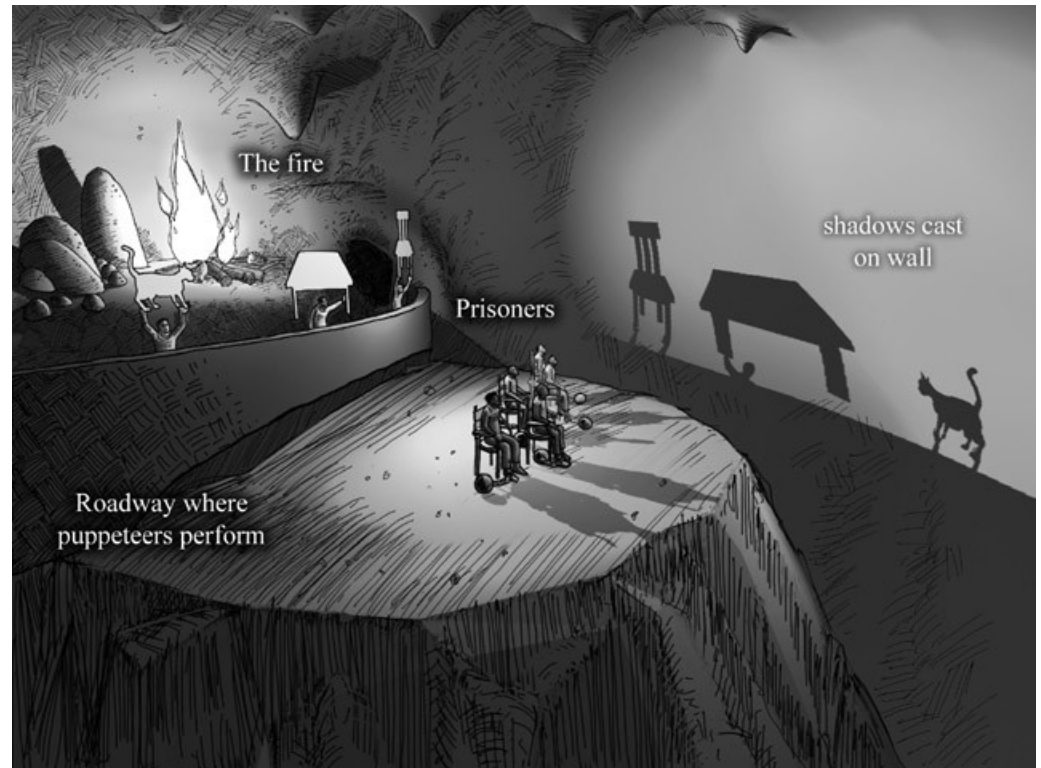


Jak vznikají informace
Rozložení dat

Význam statistické analýzy dat



- Sběr a vyhodnocování dat je způsobem k uchopení a pochopení reality.
- Chápání reality je vždy nedokonalé a nepřesné.
- Statistika umožňuje vnést do pochopení reality určitou spolehlivost a ukázat, jak je velká.



Význam statistické analýzy dat



- Realita je variabilní a statistika je věda zabývající se variabilitou.
- Korektní analýza variability a její pochopení přináší užitečné informace o realitě.
- V případě deterministického světa by statistická analýza nebyla potřebná.
- V případě zcela chaotického světa by nebyla možná.

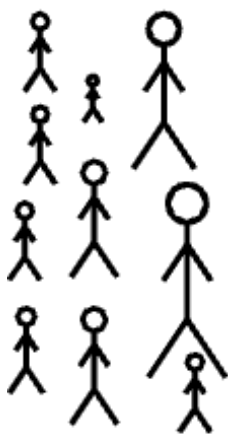


Práce s variabilitou v analýze dat

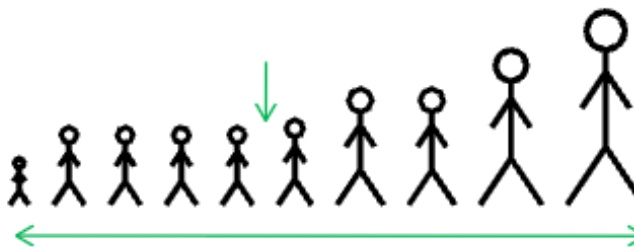


- Dva hlavní přístupy k variabilitě:

Variabilita
dat



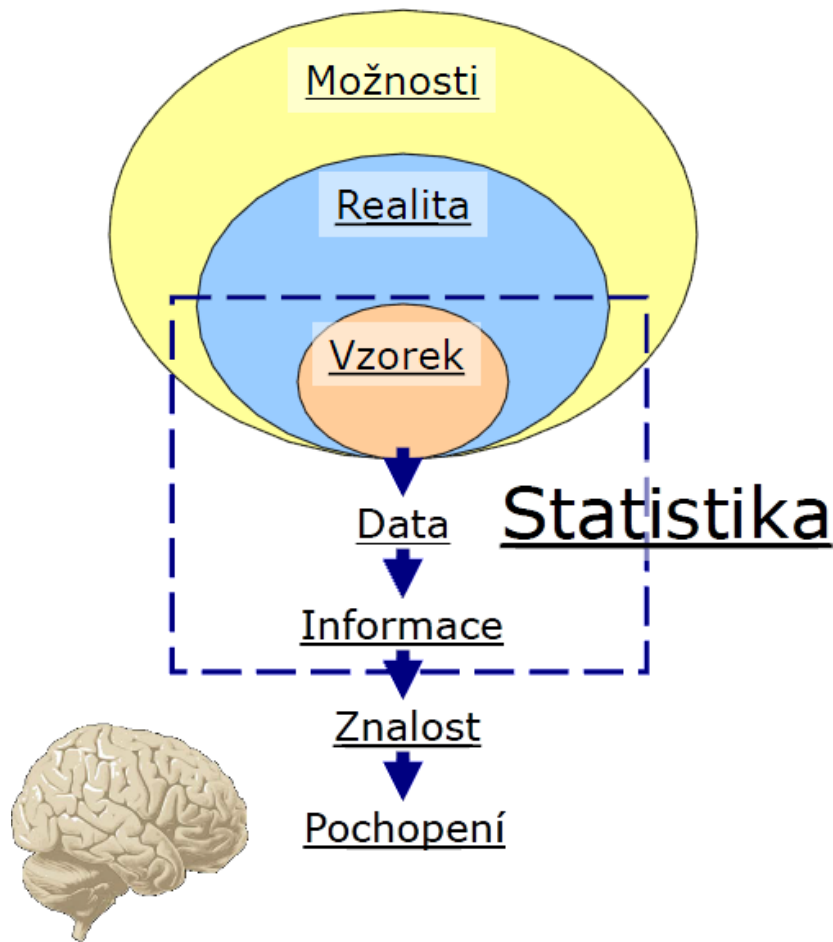
Popisná analýza: charakterizace variability



Testování hypotéz: vysvětlení variability

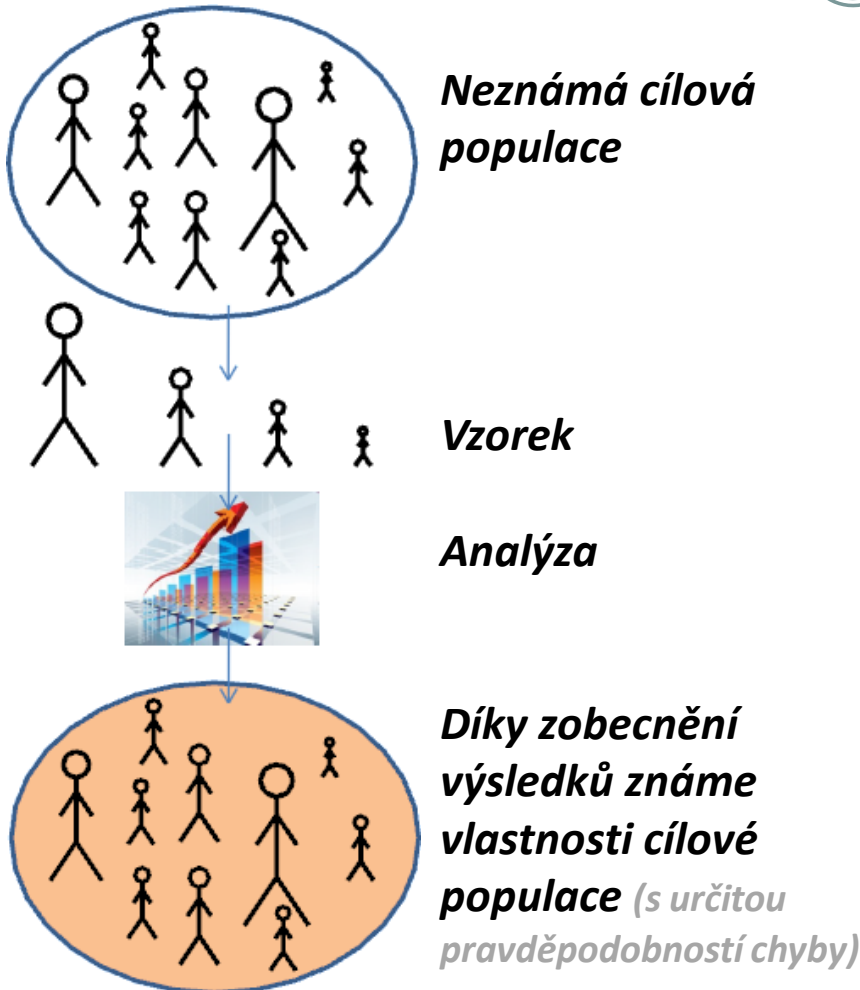


Práce s variabilitou v analýze dat



- Statistika není schopna činit závěry o jevech neobsažených ve zkoumaném vzorku.
- Statistika je nasazena v procesu získání informací ze vzorkovaných dat a je podporou v získání znalosti a pochopení problému.
- Statistika není náhradou naší inteligence!

Práce s variabilitou v analýze dat

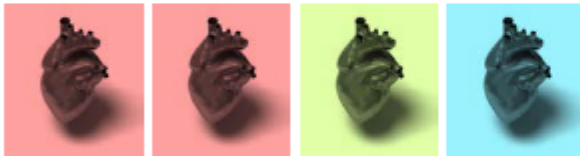


- Cílem analýzy není pouhý popis a analýza vzorku, ale zobecnění výsledků ze vzorku na jeho cílovou populaci.
- Pokud vzorek nereprezentuje cílovou populaci, vede zobecnění k chybným závěrům.

Význam vzorkování ve statistice



- Statistika hovoří o realitě prostřednictvím vzorku!!!
- Statistické předpoklady korektního vzorkování je nutné dodržet
- Náhodný výběr z cílové populace
- Representativnost: struktura vzorku musí maximálně reflektovat realitu



- Nezávislost: několikanásobné vzorkování téhož objektu nepřináší ze statistického hlediska žádnou novou informaci



Velikost vzorku a přesnost statistických výstupů



- Existuje skutečné rozložení a skutečný průměr měřené proměnné

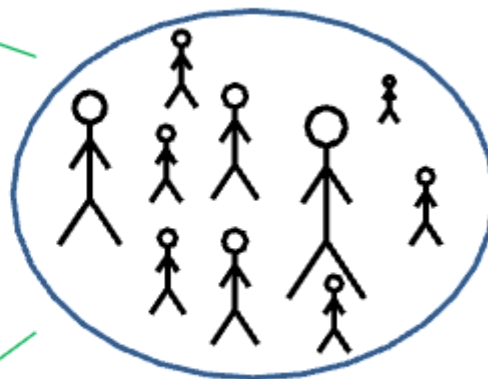
- Z jednoho měření nezjistíme nic



- Vzorek určité velikosti poskytuje odhad reálné hodnoty s definovanou spolehlivostí



- Vzorkování všech existujících objektů poskytne skutečnou hodnotu dané popisné statistiky, nicméně tento přístup je ve většině případů nereálný.



Obecné schéma aplikace statistické analýzy



Experimentální design

Jak velký vzorek je nezbytný pro statisticky relevantní výsledky?
Klíčová stratifikační kritéria cílové populace.

Vzorkování

Vzorkovací plán zabezpečující náhodnost a reprezentativnost vzorku.

Uložení a management dat

Uložení dat ve vhodné formě a jejich vyčištění předcházející vlastní analýze je klíčovým krokem statistické analýzy.

Vizualizace dat

Grafická inspekce dat je nezbytným krokem analýzy vzhledem ke schopnosti lidského mozku primárně akceptovat obrazová data. Poskytne vhled do dat, představu o jejich rozložení, vazbách proměnných apod.

Popisná analýza

Popisná analýza umožňuje vyhodnotit srovnáním s existující literaturou realističnost naměřených rozsahů dat.

Testování hypotéz

Testování vazeb mezi různými proměnnými s cílem navzájem vysvětlit jejich variabilitu a tím přispět k pochopení řešeného problému.

Modelování

Možným vyvrcholením analýzy je využití získaných znalostí a pochopení problému k vytvoření prediktivních modelů.

1a. Teoretické pozadí statistické analýzy

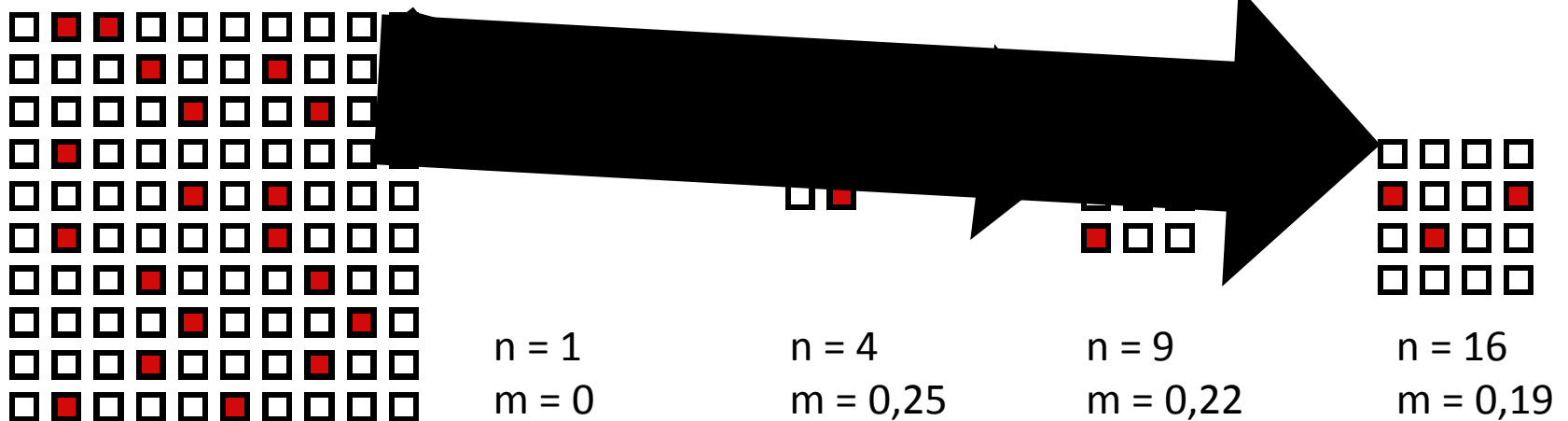


Jak vznikají informace
Rozložení dat

Anotace



- Základním principem statistiky je pravděpodobnost výskytu nějaké události. Prostřednictvím vzorkování se snažíme odhadnout skutečnou pravděpodobnost události.
- Klíčovou otázkou je velikost vzorku, čím větší vzorek, tím větší šance na projevení se skutečné pravděpodobnosti výskytu jevu (a tím je také nákladnější analýza).



Definice



Náhodný jev značíme velkým latinským písmenem, např. A . Jde o jev, pro který požadujeme tzv. statistickou stabilitu, tj. aby při n opakování pokusu platilo pro relativní četnost výsledku:

$$\lim_{n \rightarrow \infty} f(A) = \lim_{n \rightarrow \infty} \frac{n_A}{n} = p(A)$$

Elementární jev nejjemnější možný náhodný jev, tj. náhodný jev, který nelze vyjádřit jako sjednocení dvou jiných neprázdných náhodných jevů. Značí se obvykle ω .

Prostor elementárních jevů značíme obvykle Ω , jde o libovolnou neprázdnou množinu (její prvky nazýváme elementárními jevy).

Platí tedy, že elementární jevy jsou prvky prostoru elementárních jevů, rovněž jsou prvky náhodných jevů a náhodné jevy jsou podmnožiny prostoru elementárních jevů.

Definice



Ω – prostor
elementárních
jevů

A – náhodný jev

ω – elementární jev

ω – elementární jev

A – náhodný jev

ω – elementární jev

A – náhodný jev

ω – elementární jev

Definice



σ -algebra systém (množina) podmnožin prostoru elementárních jevů A (označujeme \mathcal{A}) splňující následující podmínky:

1. \mathcal{A} je neprázdná množina,
2. $A \in \mathcal{A} \Rightarrow \mathcal{A} \setminus A \in \mathcal{A}$
3. sjednocení libovolného počtu $A_i \in \mathcal{A}$.

Jevové pole uspořádaná dvojice prostoru elementárních jevů a na něm definované σ -algebry (Ω, \mathcal{A}) . Jevové pole se také někdy nazývá měřitelný prostor.

Pravděpodobnost reálná množinová funkce P definovaná na množině A σ -algebry (Ω, \mathcal{A}) tak, že jsou dodrženy následující podmínky:

(podle Kolmogorova)

1. $P(\Omega) = 1$
2. $\forall A \in \mathcal{A}: P(A) \geq 0$
3. pravděpodobnost součtu neslučitelných jevů je rovna součtu pravděpodobnosti těchto neslučitelných jevů.

Definice



- Pravděpodobnostní prostor** uspořádaná trojice prostoru elementárních jevů, na něm definované σ -algebry a jim příslušné pravděpodobnostní funkce (Ω, \mathcal{A}, P) .
- Borelovská σ -algebra** je σ -algebra \mathcal{B} generovaná systémem borelovských množin S , tj. množin splňujících podmínku:
1. $S = (-\infty, x)$, kde $x \in \mathbb{R}$.
- Náhodná veličina** reálná množinová funkce X definovaná na prostoru elementárních jevů Ω nějakého pravděpodobnostního prostoru (Ω, \mathcal{A}, P) , splňující pro nějakou borelovskou σ -algebru \mathcal{B} předpoklad:
1. $B \in \mathcal{B} \Rightarrow \{\omega \in \Omega: X(\omega) \in B\} \in \mathcal{A}$.

Pravděpodobnostní prostor je měřitelný prostor s přidanou funkcí pravděpodobnosti.

Definice



Náhodná veličina se někdy také nazývá náhodná proměnná nebo měřitelná funkce, borelovské množiny se někdy též nazývají měřitelné množiny.

Lze ukázat, že dostatečnou podmínkou pro to, aby X byla náhodná veličina je vztah $\forall x \in \mathbb{R}: \{X < x\} \in \mathcal{A}$.

Rozdělení pravděpodobnosti

množinová funkce, která každé borelovské množině B přiřadí pravděpodobnost tak, že je dodržena následující podmínka:

1. $P_X(B) = P(\{\omega \in \Omega: X(\omega) \in B\})$ pro $B \in \mathcal{B}$.

Náhodná veličina přiřazuje náhodným jevům měřitelné hodnoty (reálná čísla), rozdělení pravděpodobnosti pak každé takové hodnotě (reprezentované nějakou borelovskou množinou B) přiřazuje pravděpodobnost, tj. hodnotu mezi 0 a 1 takovou, že jsou dodrženy předpoklady po definici pravděpodobnosti uvedené dříve.

Definice



Ω – prostor
elementárních
jevů

Jevové pole

\mathcal{A} – množinová
 σ -algebra

A – náhodný jev

ω – elementární jev

ω – elementární jev

\mathcal{B} – borelovská
 σ -algebra

1

P – pravděpodobnost

X – náhodná veličina

P_x – rozdělení pravděpodobnosti

0

$-\infty$

B – borelovské množiny

JAK vznikají informace ? základní pojmy

Skutečnost

Náhoda
(vybere jednu z možností pokusu)

Jev

podmnožina množiny všech možných výsledků
(elementárních jevů) pokusu/děje, o které lze
říct, zda nastala nebo ne

Pozorovatel

Rozliší, co nastalo

- a) podle možností
- b) podle toho, jak potřebuje

Jevové pole

třída všech jevů, které jsme se rozhodli nebo
jsme schopni sledovat

Skutečnost + Jevové pole = Měřitelný prostor

Experimentální jednotka - objekt, na kterém se provádí šetření

Populace - soubor experimentálních jednotek Znak - vlastnost sledovaná na objektu

Sledovaná veličina - číselná hodnota vyjadřující výsledek náhodného experimentu

Znak se stává náhodnou veličinou, pokud se jeho hodnota zjišťuje
vylosováním objektu ze základního souboru

Výběr - výběrová populace - cílová populace

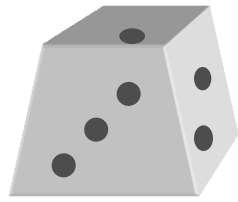
Náhodný výběr

Reprezentativnost

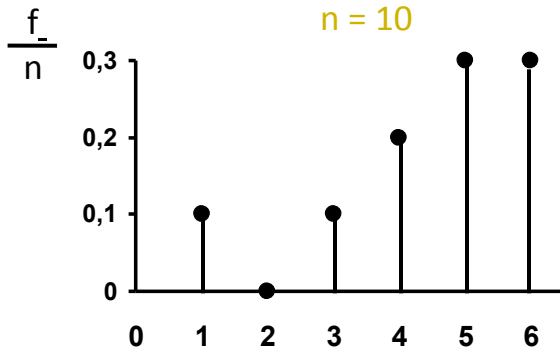
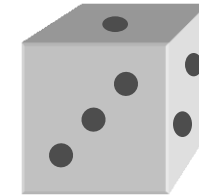
JAK vznikají informace ?

„Empirical approach“

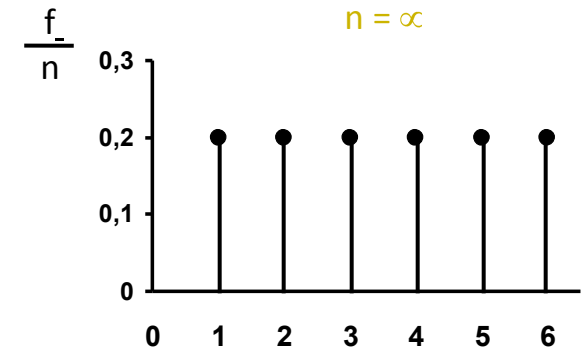
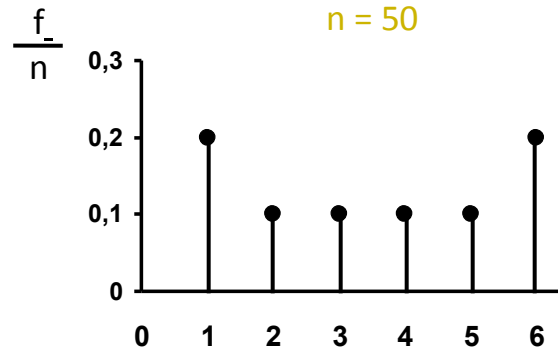
„Classical approach“



Empirický postup



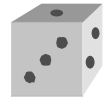
možné jevy: čísla 1 – 6



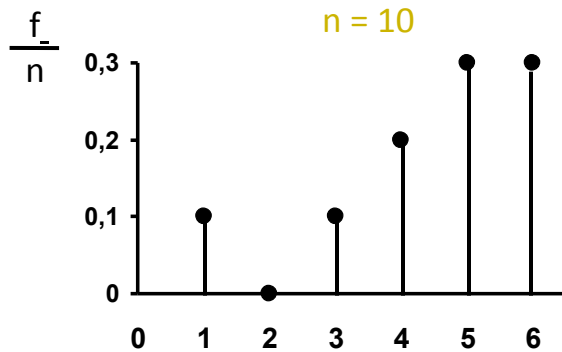
n – počet hodů (opakování)

U složitých stochastických systémů se pravdě blížíme až po odvedení značného množství experimentální práce: musíme dát systému šanci se projevit

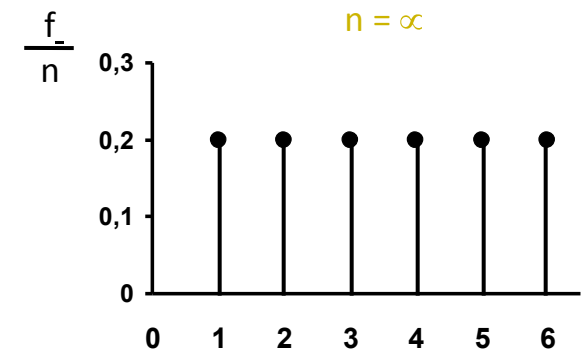
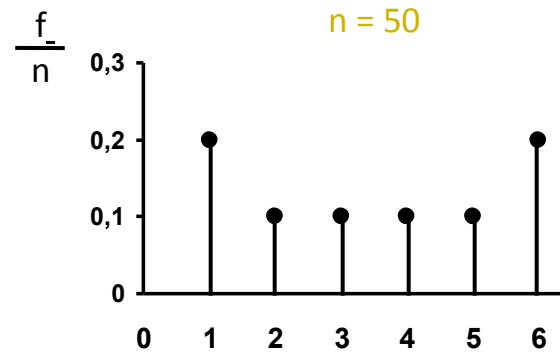
JAK vznikají informace ?



Empirický postup



možné jevy: čísla 1 – 6



n – počet hodů (opakování)



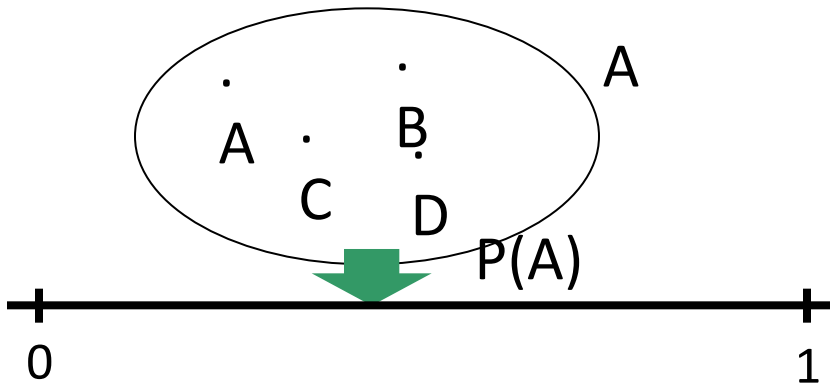
Při realizaci náhodného experimentu roste se zvyšujícím se počtem opakování pravdivá znalost systému (výsledky se stávají stabilnější) ...diskutabilní je ale ovšem míra zobecnění konkrétního experimentu

Empirický zákon velkých čísel



Při opětovné nezávislé realizaci téhož náhodného experimentu se podíl výskytů sledovaného jevu mezi všemi dosud provedenými realizacemi zpravidla ustaluje kolem konstanty.

Pravděpodobnost je libovolná reálná funkce definovaná na jevovém poli A , která každému jevu A přiřadí nezáporné reálné číslo $P(A)$ z intervalu $0 - 1$.



Z praktického hlediska je
pravděpodobnost
idealizovaná relativní četnost

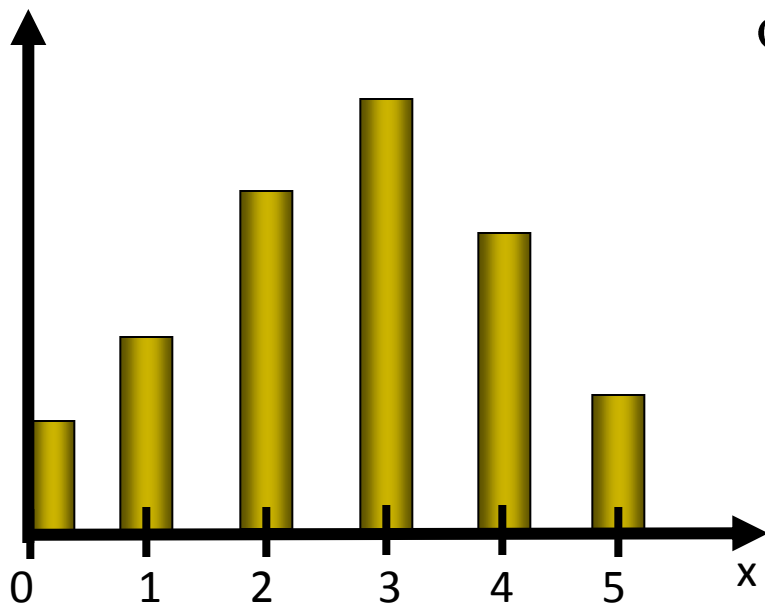
- $P(A) = 1$ jev jistý
- $P(A) = 0$ jev nemožný
- $P(A \cap B) = P(A) \cdot P(B)$ nezávislé jevy
- $P(A \cap B) = P(A) \cdot P(B/A)$ závislé jevy
- $P(A/B) = P(A \cap B) / P(B)$ podmíněná pravděpodobnost

Pravděpodobnost výskytu jevu – rozložení dat



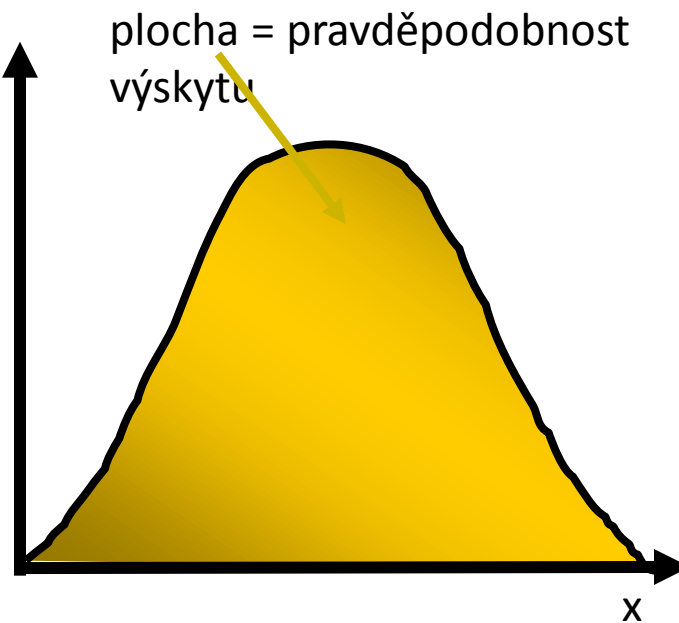
- ✦ existuje pravděpodobnost výskytu jevů (nedeterministické závěry)
- ✦ „vše je možné“: pouze jev s pravděpodobností 0 nikdy nenastane
- ✦ pravděpodobnost lze zkoumat retrospektivně i prospektivně

pravděpodobnost
výskytu



počet chlapců v rodině s X dětmi

$\varphi(x)$

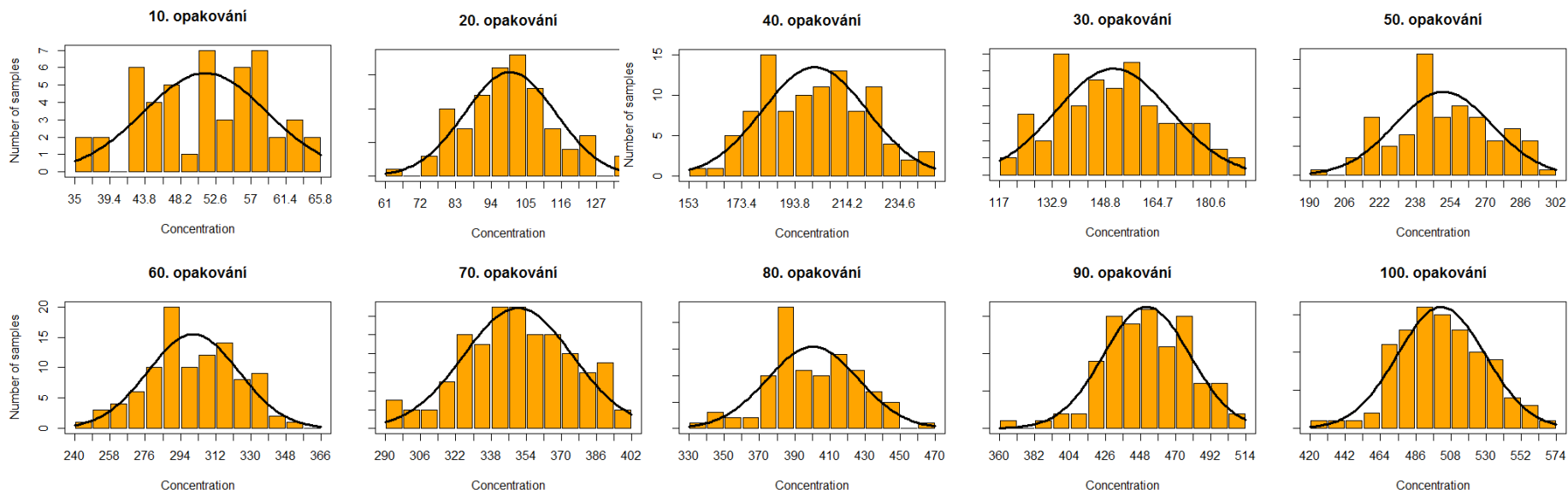


výška postavy

Centrální limitní věta



Pokud lze náhodnou veličinu X vyjádřit jako součet náhodných veličin X_1, X_2, \dots, X_n , které mají shodné rozdělení, konečnou střední hodnotu a konečný rozptyl, platí, že rozdělení veličiny X se vzrůstajícím n konverguje (poměrně rychle) k normálnímu rozdělení.



Instalace R



Webová stránka <https://cran.r-project.org/>



The Comprehensive R Archive Network

CRAN

[Mirrors](#)

[What's new?](#)

[Task Views](#)

[Search](#)

About R

[R Homepage](#)

[The R Journal](#)

Software

[R Sources](#)

[R Binaries](#)

[Packages](#)

[Other](#)

Documentation

[Manuals](#)

[FAQs](#)

[Contributed](#)

Download and Install R

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- [Download R for Linux](#)
- [Download R for \(Mac\) OS X](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

Source Code for all Platforms

Windows and Mac users most likely want to download the precompiled binaries listed in the upper box, not the source code. The sources have to be compiled before you can use them. If you do not know what this means, you probably do not want to do it!

- The latest release (Tuesday 2016-06-21, Bug in Your Hair) [R-3.3.1.tar.gz](#), read [what's new](#) in the latest version.
- Sources of [R alpha and beta releases](#) (daily snapshots, created only in time periods before a planned release).

Instalace R studia



Webová stránka

<https://www.rstudio.com/products/rstudio/download3/>

The screenshot displays the RStudio environment with the following components:

- Source Editor:** Contains R code for loading the 'diamonds' dataset, summarizing it, and creating a faceted scatter plot 'p' of Price vs. Carat, faceted by Clarity.
- Console:** Shows the execution output, including summary statistics for the 'diamonds' dataset and the 'price' variable.
- Workspace:** Lists the loaded data object 'diamonds' (53940 observations) and the plot object 'p'.
- Plots Panel:** Displays the 'Diamond Pricing' scatter plot, where points are colored by Clarity (I1, SI2, SI1, VS2, VS1, VVS2, VVS1, IF).

```
1 library(ggplot2)
2 source("plots/FormatPlot.R")
3
4 View(diamonds)
5 summary(diamonds)
6
7 summary(diamonds$price)
8 aveSize <- round(mean(diamonds$carat), 4)
9 clarity <- levels(diamonds$clarity)
10
11 p <- qplot(carat, price,
12            data=diamonds, color=clarity,
13            xlab="Carat", ylab="Price",
14            main="Diamond Pricing")
15
```

x	y	z
Min. : 0.000	Min. : 0.000	Min. : 0.000
1st Qu.: 4.710	1st Qu.: 4.720	1st Qu.: 2.910
Median : 5.700	Median : 5.710	Median : 3.530
Mean : 5.731	Mean : 5.735	Mean : 3.539
3rd Qu.: 6.540	3rd Qu.: 6.540	3rd Qu.: 4.040
Max. : 10.740	Max. : 58.900	Max. : 31.800

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
326	950	2401	3933	5324	18820