

Jazyk R

II. Dáta, Grafy, Balíčky, Štatistika



Ako Začneme?

1. Registrácia na Jetbrains Datalore +

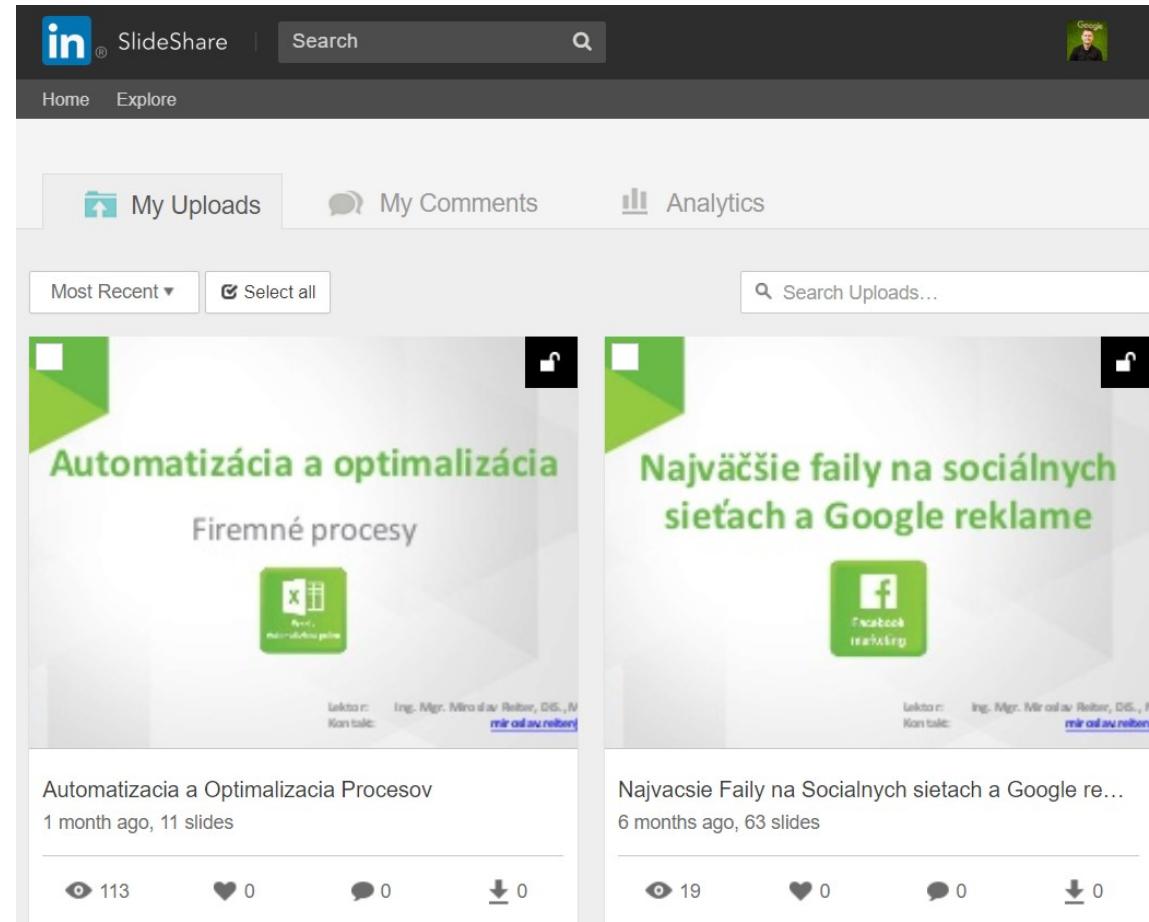
Stiahni si R a Anaconda Navigator

- <https://datalore.jetbrains.com/>
- <https://cran.r-project.org/bin/windows/base/>
- <https://www.anaconda.com/products/>

2. Pridajte si ma na LinkedIn

- www.linkedin.com/in/miroslav-reiter

3. Prezentácia a materiály po prednáške



[Home](#) » Štatistika v R - Spracovanie a vizualizácia dát

Štatistika v R - Spracovanie a vizualizácia dát

Kurz sa zameriava na pokročilejšie techniky práce s prostredím pre štatistickú analýzu R. Nosnou tému kurzu bude spracovanie a vizualizácia dát.

1. Spracovanie a vizualizácia dát v R:

- CSV a Excel súbory
- Koláčové grafy
- Stípcové grafy
- Krabicové grafy
- Histogramy
- Čiarové grafy
- Bodové grafy

2. Dátové štruktúry a programovanie

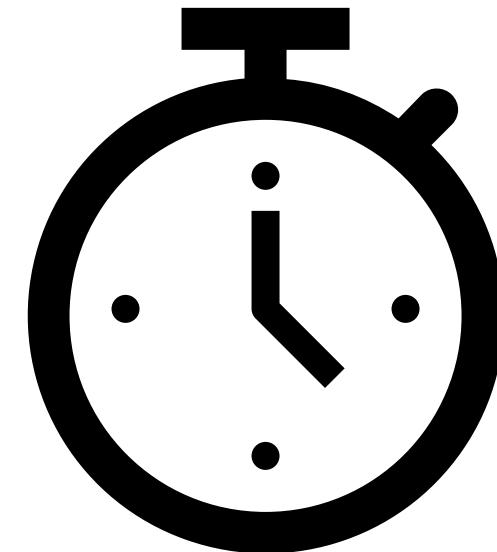
- Matice a Submatice
- Polia
- Faktory
- Dátové rámce/Tabuľky
- Balíčky
- GUI

3. Štatistická analýza

- Korelačná analýza
- Regresná analýza
- Analýza kontingenčných tabuľiek
- Časové rady

Úvodné informácie

- Časový rozvrh (9:00-13:30)
 - Prestávky
 - Mobilné telefóny a zariadenia
-
- Priprav si otázky a rovno sa pýtaj
 - Interaktívna forma



O lektorovi - Miroslav Reiter

10000+
klientov a
500+ firiem

Programátor
Analytik
Manažér

Google
Microsoft
ISTQB tréner

115
certifikácií

83 príručiek a
publikácií

13 škôl

50+
projektov

Vlastná firma



MOTIVÁCIA

Študuje 5 odborov a absolvoval už 12 univerzít. Ako zvláda stres a manažovanie času?



Foto: Jakub Kovalík pre FMK UCM | Miroslav Reiter na prednáške Grow with Google na FMK UCM.

**Nikola Kotláriková**

19. júl 2022 · 8 min. čítania



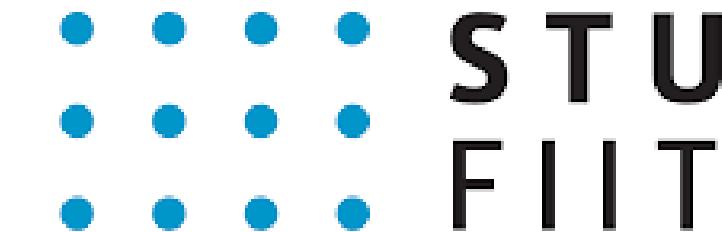
Miroslav Reiter



1. PhDr. VŠM (Podnikovný manažment)
2. Ing. STU FEI (**Aplikovaná informatika**)
3. Mgr. UK FM (**Strategický manažment a marketing**)
4. Mgr. VŠM (**Manažérstvo kvality**)
5. Mgr. VŠEMVŠ (Verejná správa)
6. Mgr. DTI (Učiteľstvo ekonomických predmetov)
7. DiS. AMOS (Cestovný ruch)
8. MBA LIGS (Executive management)
9. DBA Humanum (**IT manažment**)
10. MPA IES (Verejná správa a samospráva krajov)
11. MSc. Humanum (**Bezpečnosť informačných systémov**)
12. Ing. Paed. IGIP
13. Mgr. PEVŠ (**Bezpečnosť informačných systémov**)



DIGITÁLNA
UNIVERZITA



FAKULTA MANAGEMENTU
Univerzity Komenského
v Bratislave



Vyberte si online kurz

Naučte sa programovať, tvoriť webstránky a grafiku, manažovať alebo sa zamerajte na osobný rozvoj. Všetko jednoducho vďaka našim online kurzom z pohodlia domova.

Ročné predplatné na všetky online kurzy

2299.99€

399.99€

Prístup pre Vás do všetkých aktuálnych aj pripravovaných online kurzov

12 mesačná platnosť

Kúpiť teraz



407 kurzov v ponuke



Zábavné online lekcie



Akreditované kurzy



11 rokov skúseností



Certifikovaní profesionálni lektori

Odporučame Kurzy špeciálne pre vás



Online kurz SAP I.
Začiatočník
224,00€ 292,50€



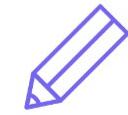
Online kurz Java I.
Začiatočník
67,00€ 88,40€



Online kurz PRINCE2 Foundation
224,00€ 308,70€



Online kurz Lektor (Akreditovaný Kurz Lektor)
127,00€ 167,50€



Online kurz Copywriting I.
Začiatočník
50,00€ 66,30€



Online kurz Testovanie Softvéru I. Začiatočník
140,00€ 183,30€

Moje začiatky

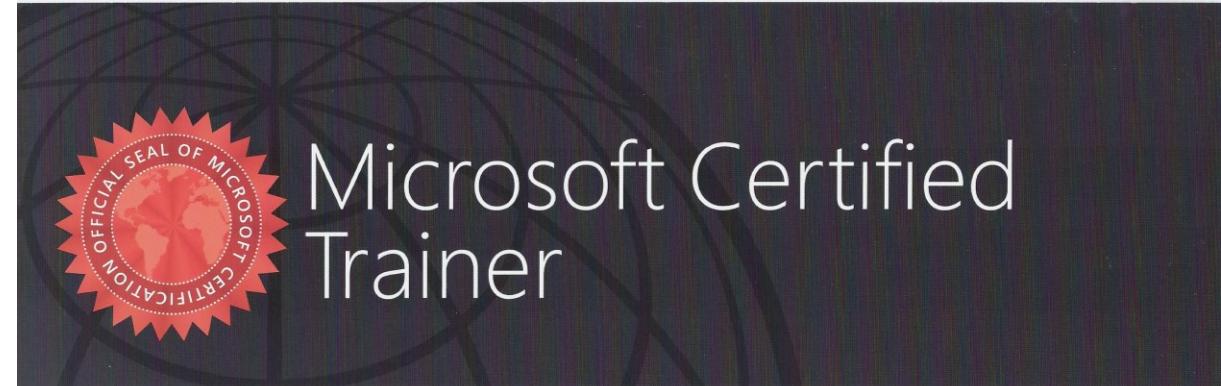


Miroslav Reiter

získava status
Google Certified Trainer

Automation

Google



MIROSLAV REITER

Has successfully completed the requirements to be recognized as a Microsoft Certified Trainer

N. S. [Signature]
Satya Nadella
Chief Executive Officer

Microsoft
CERTIFIED
Trainer



Luigi, Mário
a Yoshi



Čo robíte?

1. Študent/učiteľ

2. Zamestnanec

3. Podnikateľ

4. Nezamestnaný/materská

5. Dievča pre všetko



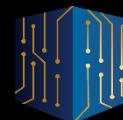
National competence centre for high performance computing
SLOVAKIA



EURO



EuroHPC
Joint Undertaking



NATIONAL
SUPERCOMPUTING
CENTRE



Vzdelávanie

Kurzy:
itkurzy.sav.sk



Propagácia

Prednášky:
<https://eurocc.nscc.sk/vzdelavanie/prednasky/prednasky-archiv/>



HPC služby

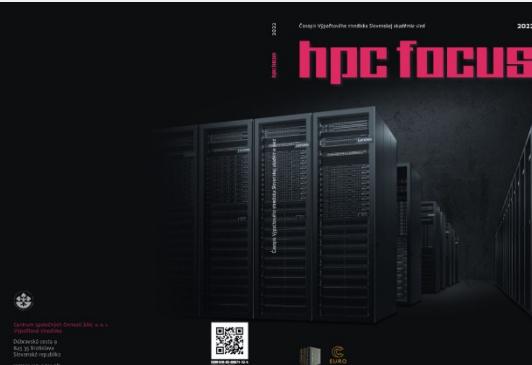
Prístup k
výpočtovým
prostriedkom



Spolupráca

Pilotné projekty

Dlhodobá
spolupráca



robíme it



Qubit
Conference

S kým spolupracujeme:

- Akademické inštitúcie, univerzity,
ústavy SAV...
- Verejná správa
- Súkromné firmy
- Tretí sektor



Sledujte nás na sociálnych sietiach:
#nccprehpc



Interaktívna prednáška

Aktívne používanie a zapájanie sa

Participants (20)

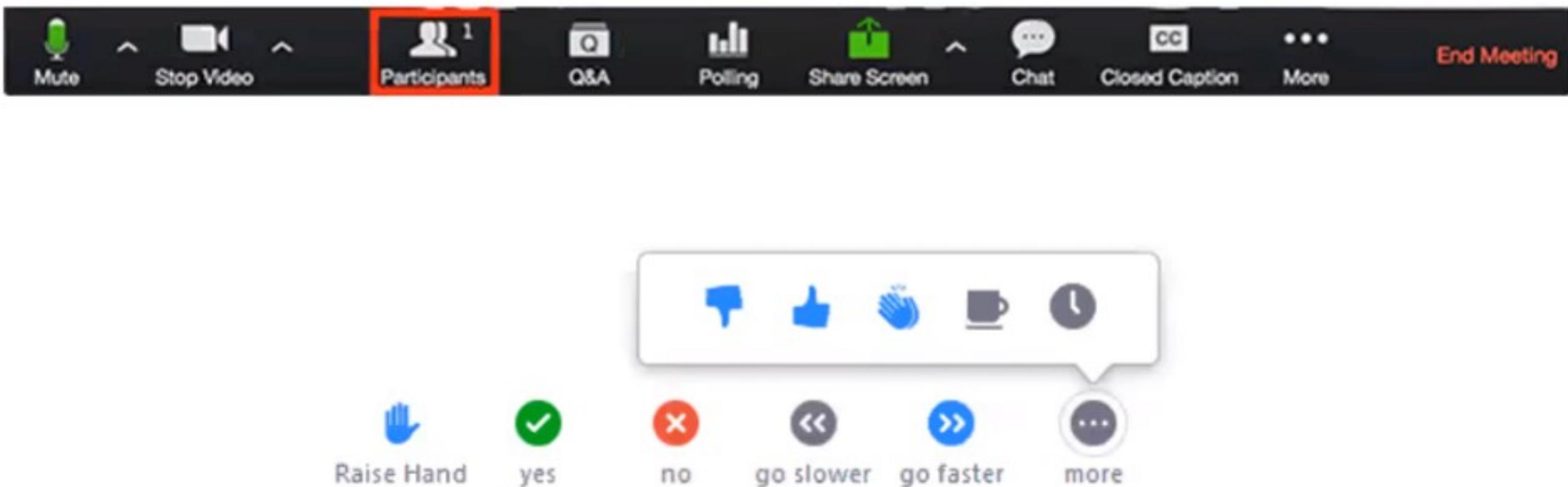
Find a participant

Participant	Microphone	Video	Hand Raised	Feedback	More
Miroslav Reiter (Me)	✓	🔇	🕒	🕒	...
	✗	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...
	▢	▢	🕒	🕒	...

Raise Hand yes no go slower go faster more

Invite Mute Me

Používame Zoom



Vaše ciele a očakávania

1. Základy jazyka R a Data Science

2. Základy analytického/štatistického myslenia

3. Základy programovania v R

4. Základy s knižnicami/balíčkami

5. Základy práce s vývojovými nástrojmi pre jazyk R

Zábava je v zaručená v každom bode :-)





Data Science a Jazyk R

Začiatky v R a Data Science



R a Data Science je ako vzťah...

[Create](#)[Home](#)[Competitions](#)[Datasets](#)[Code](#)[Discussions](#)[Learn](#)[More](#)[Your Work](#)[RECENTLY VIEWED](#)[House Prices - Advanc...](#)[Educational institutes ...](#)[Taylor Swift Spotify Data](#)[Global pollution by cou...](#)[Titanic - Ensemble Cla...](#)[RECENTLY EDITED](#)[View Active Events](#)

ZVR_842 · UPDATED 9 DAYS AGO



26

[New Notebook](#) [Download \(15 kB\)](#)

Global pollution by counties

Worldwide waste by countries (inc. all types of pollution by 2010-2020)

[Data](#) [Code \(1\)](#) [Discussion \(0\)](#)

About Dataset

Usability ⓘ

9.71

License

CC BY-SA 4.0

Expected update frequency

Annually

[Classification](#) [Data Analytics](#) [Exploratory Data Analysis](#) [Energy](#) [Advanced](#) [Global](#)**country_level_data_0.csv** (46.19 kB)[Detail](#) [Compact](#) [Column](#)

10 of 49 columns

About this file

there are 49 columns, be accurate with columns name because it is divided by types of wasting

region_id	country_name	A gdp	A composition_foo...	A composition_glas...	A compo...
-----------	--------------	-------	----------------------	-----------------------	------------

Data Explorer

46.19 kB

country_level_data_0.csv

Vyhľadať dataset



ditapa

NAJLEPŠÍ PROJEKT DIGITALIZÁCIE SPOLOČNOSTI 2020

Doprava



148 datasetov

Ekonomika a práca



37 datasetov

Infraštruktúra, výstavba a bývanie



46 datasetov

Kultúra, šport, cestovný ruch a Wifi



71 datasetov

Obyvateľstvo



123 datasetov

Politika a volby



156 datasetov

Priestorové údaje



82 datasetov

Rozpočet, dane a zmluvy



40 datasetov

Sociálna oblasť

Veda a vzdelávanie

Zákony a spravodlivosť

Zdravie, životné prostredie a klíma

[Welcome](#)[Preface](#)[0.1 Recipes](#)[0.2 Software and Platform Notes](#)[0.3 Conventions Used in This Book](#)[0.4 Using Code Examples](#)[0.5 How to Contact Us](#)[0.6 Acknowledgments](#)[1 R Basics](#)[1.1 Installing a Package](#)[1.2 Loading a Package](#)[1.3 Upgrading Packages](#)[1.4 Loading a Delimited Text Data...](#)[1.5 Loading Data from an Excel File](#)[1.6 Loading Data from SPSS/SAS...](#)[1.7 Chaining Functions Together ...](#)[2 Quickly Exploring Data](#)[2.1 Creating a Scatter Plot](#)[2.2 Creating a Line Graph](#)[2.3 Creating a Bar Graph](#)[2.4 Creating a Histogram](#)[2.5 Creating a Box Plot](#)[2.6 Plotting a Function Curve](#)[3 Bar Graphs](#)[3.1 Making a Basic Bar Graph](#)

R Graphics Cookbook, 2nd edition

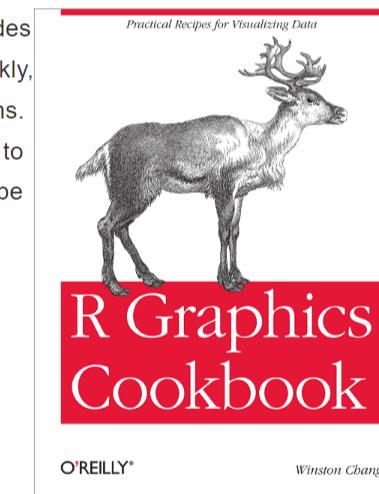
Winston Chang

2022-11-23

Welcome

Welcome to the **R Graphics Cookbook**, a practical guide that provides more than 150 recipes to help you generate high-quality graphs quickly, without having to comb through all the details of R's graphing systems. Each recipe tackles a specific problem with a solution you can apply to your own project, and includes a discussion of how and why the recipe works.

Read online here for free, or buy a physical copy on [Amazon](#).



1.1 Where did this book come from?

1.2 Who is this book for?

1.3 Why is R so great?

1.4 Why R is like a relationship...

1.5 R resources

1.6 Who am I?

1.7 Contributions and Acknowledgments

2 Getting Started

2.1 Installing Base-R and RStudio

2.2 The four RStudio Windows

2.3 Packages

2.4 Reading and writing Code

2.5 Debugging

3 Jump In!

3.1 Exploring data

3.2 Descriptive statistics

3.3 Plotting

3.4 Hypothesis tests

3.5 Regression analysis

3.6 Bayesian Statistics

3.7 Wasn't that easy?!

4 The Basics

4.1 The command-line (Console)

4.2 Writing R scripts in an editor

4.3 A brief style guide: Commenting and documentation

4.4 Objects and functions

YaRrr! The Pirate's Guide to R

Nathaniel D. Phillips

2018-01-22

Chapter 1 Preface



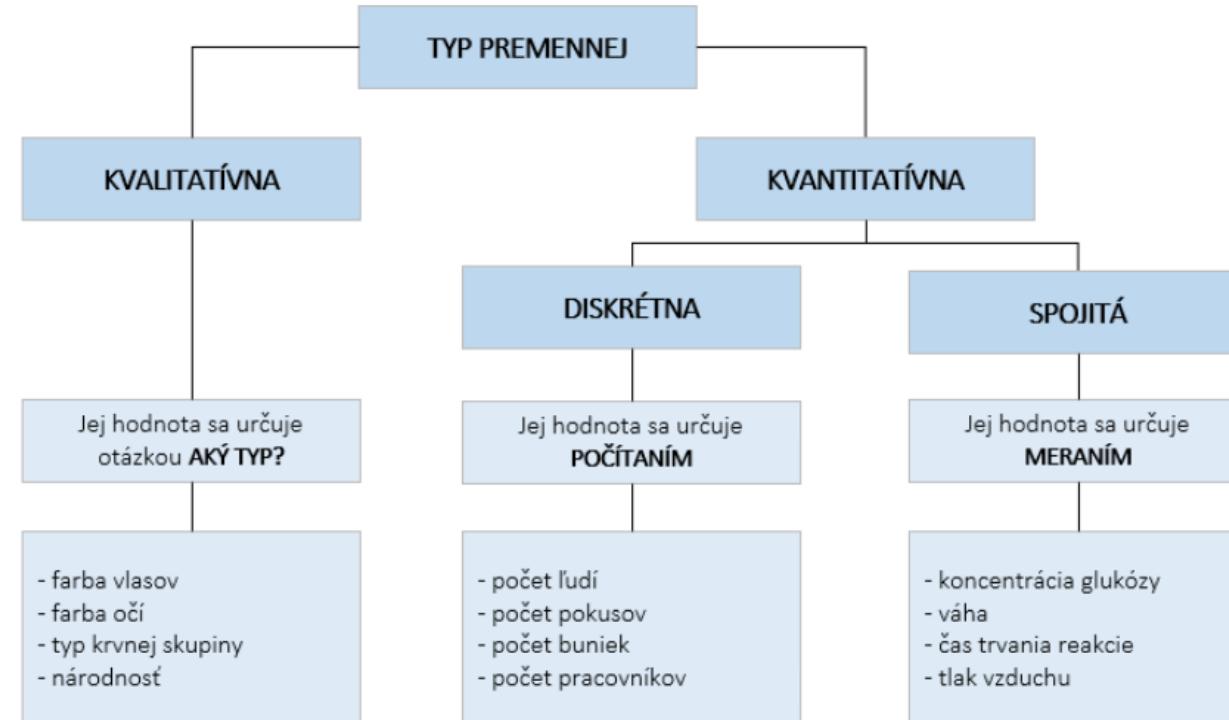
Štatistické Znaky – Premenné (Variables)

Slovné – Kvalitatívne

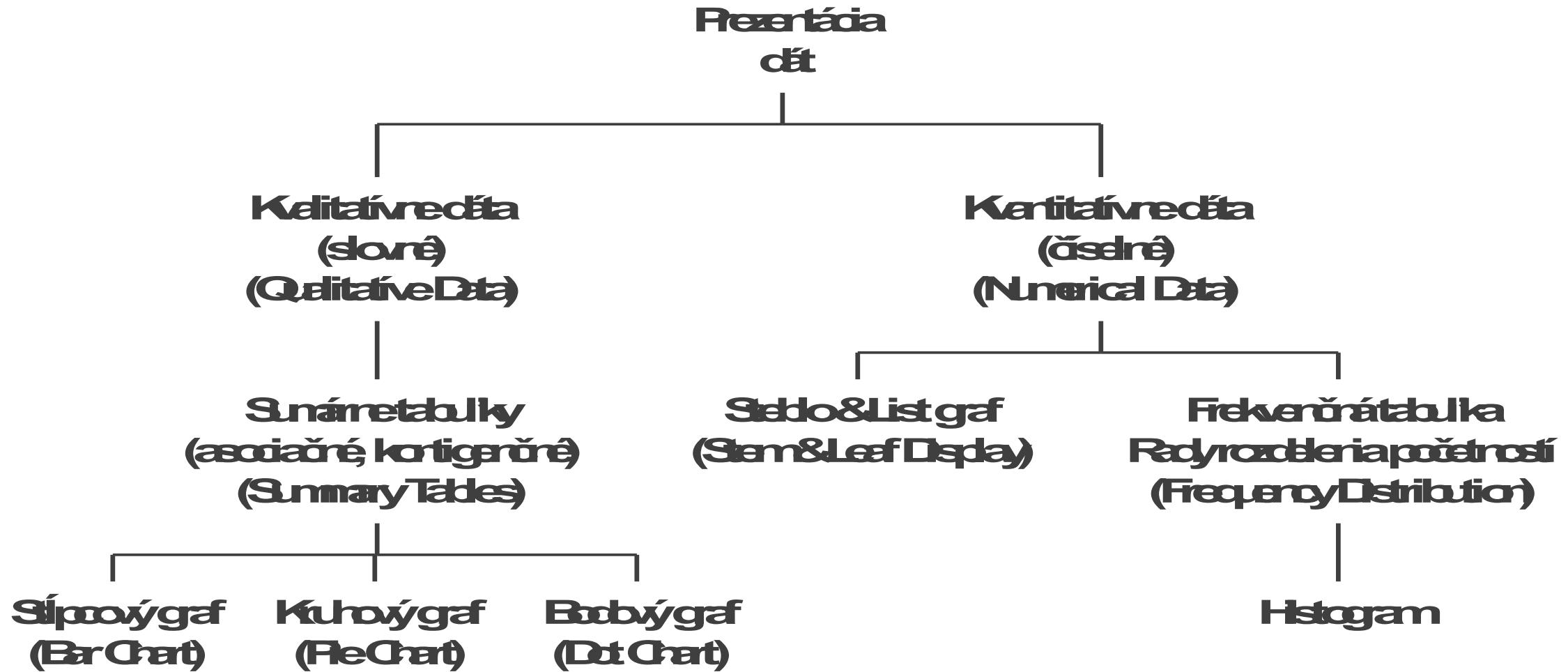
- Označenie: A, B, C, ...

Číselné – Kvantitatívne

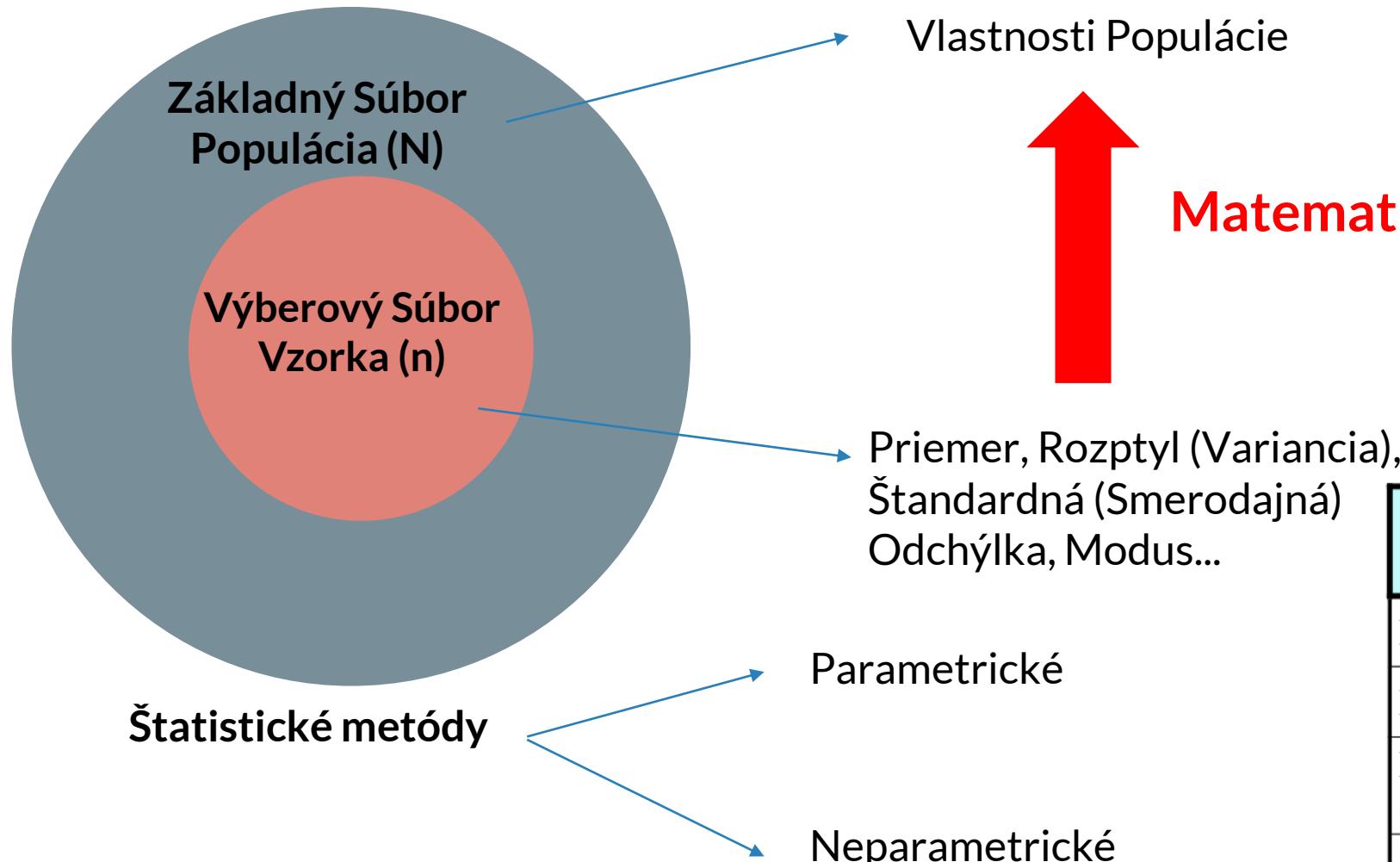
- Označenie: X, Y, Z, ...
- Diskrétné a Spojité



Spôsoby Vizualizácie (Prezentácie) Dát



Popisná (Deskriptívna) Štatistika



Ukazovateľ Measure	Populácia Population	Výber Sample
Rozsah (Size)	N	n
Priemer (Mean)	μ	\bar{x}
Rozptyl (Variance)	σ^2	s^2
Štand.oddchýlka (Stand.Deviation)	σ	s

Priemer, Medián, Modus

MEAN GIRLS

MEDIAN GIRLS

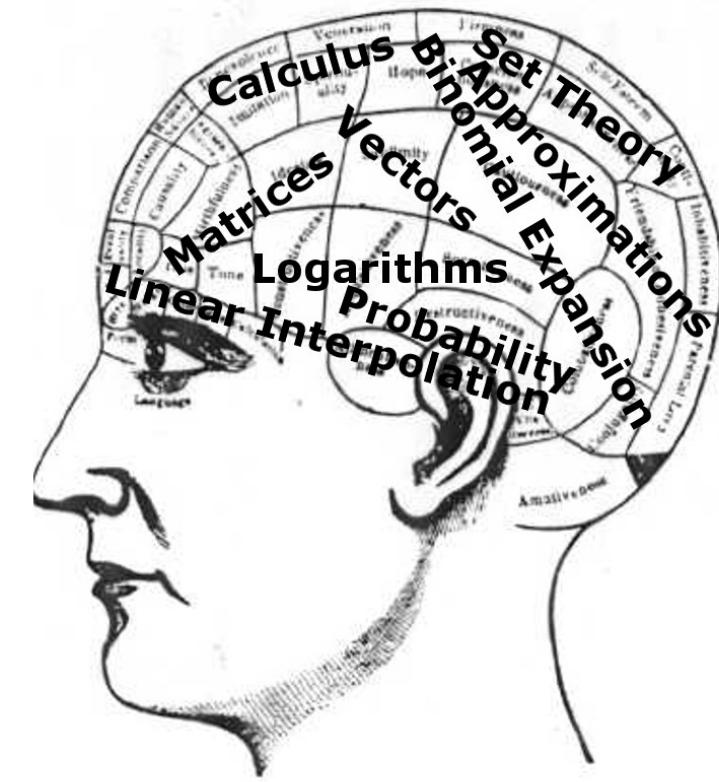
MODE GIRLS

$$\text{GIRLS}_{\text{mean}} = \frac{\sum_{i=1}^n \text{GIRLS}_i}{n}$$

$$B_m + I \cdot \left(\frac{\frac{n}{2} - (\sum \text{GIRLS}_i)_0}{\text{GIRLS}_m} \right)$$

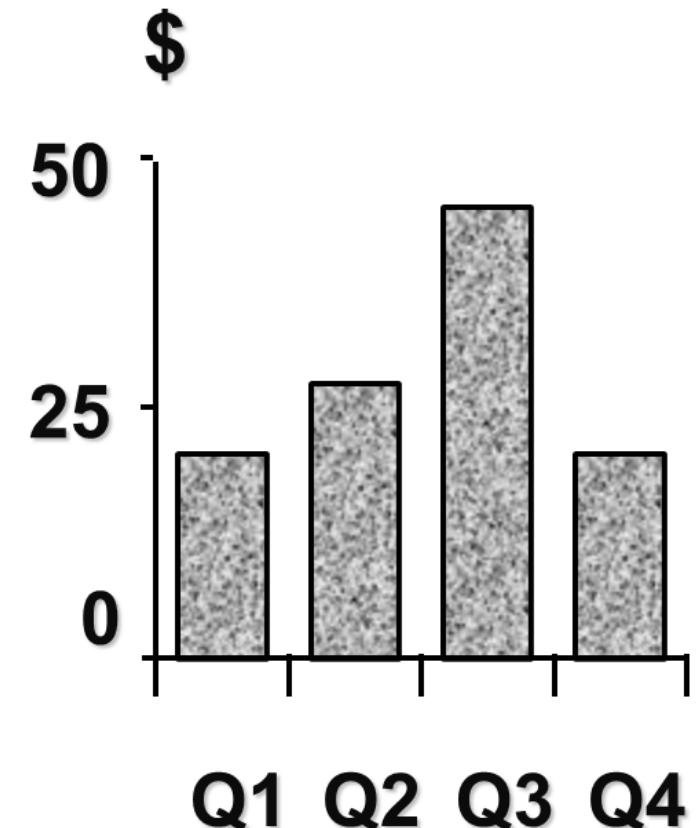
$$l + h \left(\frac{\text{GIRLS}_m - \text{GIRLS}_1}{2 \text{GIRLS}_m - \text{GIRLS}_1 - \text{GIRLS}_2} \right)$$

"Sometimes I wonder what goes on in that head of yours..."



Popisná (Deskriptívna) Štatistika

- 1. Obsah
 - Zber dát
 - Prezentácia (Vizualizácia) dát
 - Charakter dát
- 2. Ciel'
 - Popis dát
- Spracovanie Dát
 - Kontrola dát – Formálna, Vecná
 - Triedenie - usporiadanie dát do skupín (tried) podľa určitého štat. znaku (-ov) tak, aby čo najlepšie vynikli vlastnosti javu (-ov)
 - Dátové Konverzie (CSV, XLSX, XML, STATA, SAV)



$$\bar{X} = 30.5$$

$$S^2 = 113$$

MIERY

POLOHY

- Minimum
- Maximum
- Priemer
- Medián
- Modus
- Kvantity

VARIABILITY

- Štandardná odchylka
- Rozptyl
- Štandardná chyba
- Variačné rozpätie
- Medzikvartilové rozpätie
- Variačný koeficient

TVARU

- Šikmost'
- Špicatost'

MIERY

POLOHY

- Minimum
- Maximum
- Priemer
- Medián
- Modus
- Kvantity

$$\begin{array}{c} x_{\min} \\ x_{\max} \end{array}$$

$$r = \frac{n + 1}{2}$$

Q_1	Q_2	Q_3
25%	25%	25%

Pozícia i-ho kvartilu:

$$k_{Qi} = \frac{i(n + 1)}{4}$$

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

– Vážený (z frekv. tabuľky):

$$\bar{x} = \frac{1}{n} \sum_{i=1}^m x_i n_i = \frac{1}{n} (x_1 n_1 + x_2 n_2 + \dots + x_m n_m)$$

Dáta: 10, 10, 11, 13, 9, 10, 10, 8

$x_{Mo}=10$ (1 modus)

Dáta: 10, 10, 11, 12, 12, 8, 9, 10, 12

$x_{Mo}=10$ a 12 (2 modusy)

Dáta: 10, 12, 8, 9, 11, 13, 7 - dátá bez modusu

	Ab	Murder	Ab	Assault
		Min. : 0.800		Min. : 45.0
		1st Qu.: 4.075		1st Qu.: 109.0
		Median : 7.250		Median : 159.0
		Mean : 7.788		Mean : 170.8

USArrests

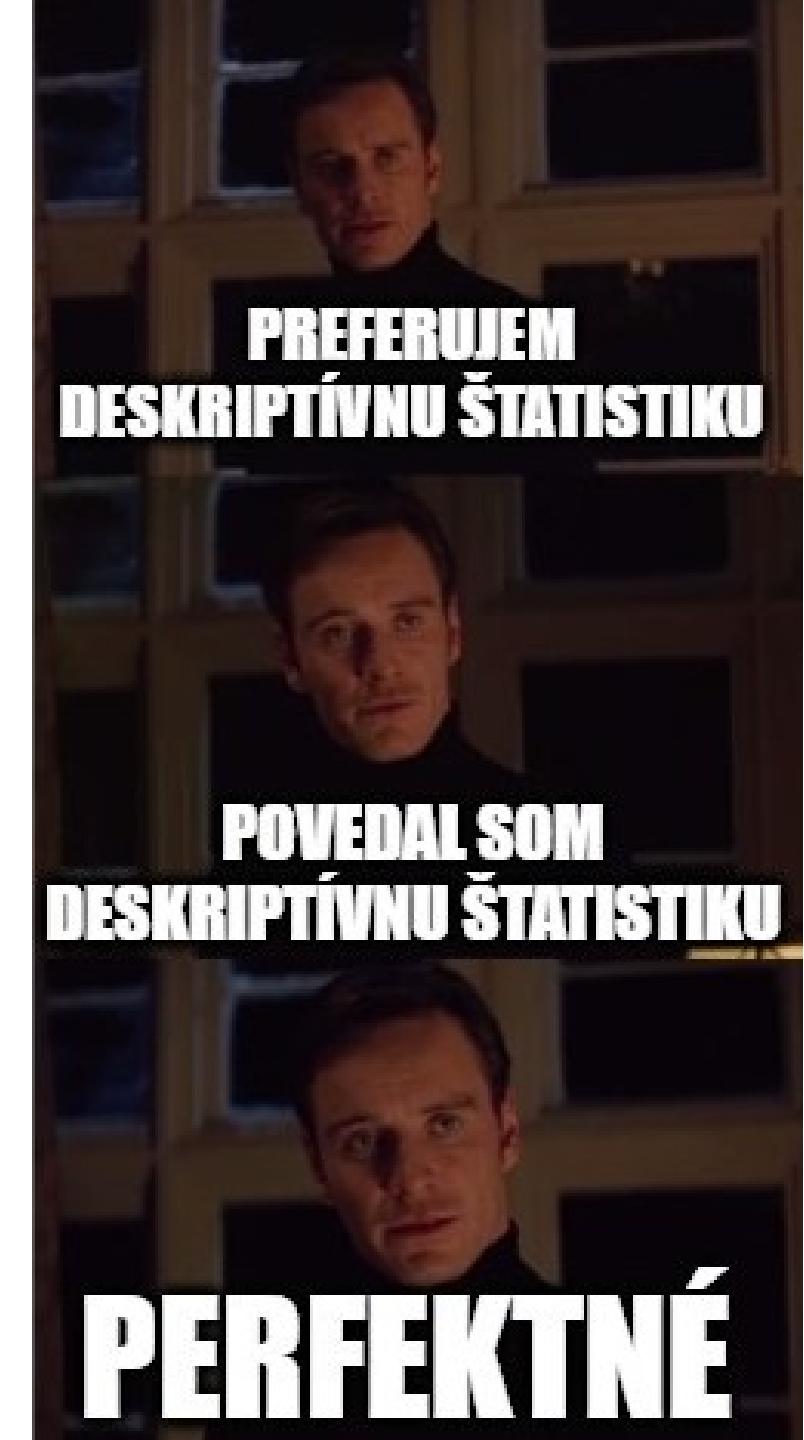
4 Variables 50 Observations

Murder

n	missing	distinct	Info	Mean	Gmd	.05	.10
50	0	43	1	7.788	5.022	2.145	2.560
.25	.50	.75	.90	.95			
4.075	7.250	11.250	13.320	15.400			

lowest : 0.8 2.1 2.2 2.6 2.7 , highest: 13.2 14.4 15.4 16.1 17.4

	3	Murder	3	Assault
nbr.val		50		50
nbr.null		0		0
nbr.na		0		0
min		0.8		45
max		17.4		337
range		16.6		292
sum		389.4		8538



Deskriptívna Štatistika

```
# get means for variables in data frame mydata  
# excluding missing values  
sapply(mydata, mean, na.rm=TRUE)
```

```
sapply(USArrests, mean)  
Murder Assault UrbanPop Rape  
7.788 170.760 65.540 21.232
```

```
# mean, median, 25th and 75th quartiles, min, max  
summary(mydata)
```

```
# Tukey min, lower-hinge, median, upper-hinge, max  
fivenum(x)
```

Deskriptívna Štatistika

```
> describe(USArrests)
```

USArrests

4 Variables 50 observations

Murder

	n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50	.75
	50	0	43	1	7.788	5.022	2.145	2.560	4.075	7.250	11.250
	.90	.95									
	13.320	15.400									

lowest : 0.8 2.1 2.2 2.6 2.7, highest: 13.2 14.4 15.4 16.1 17.4

Assault

	n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50	.75
	50	0	45	1	170.8	96.44	50.25	56.90	109.00	159.00	249.00
	.90	.95									
	279.60	297.30									

lowest : 45 46 48 53 56, highest: 285 294 300 335 337

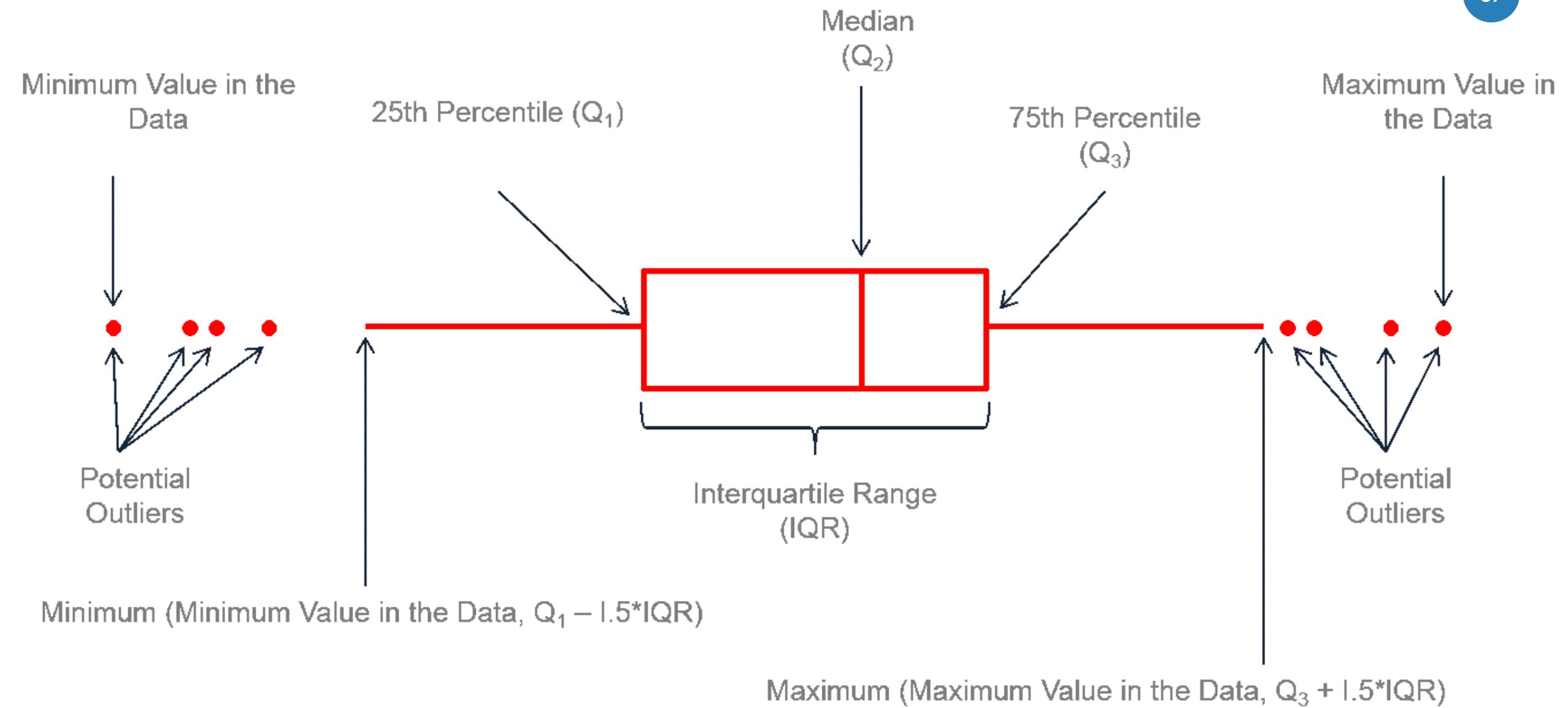
Deskriptívna Štatistika

Using the [Hmisc](#) package

```
library(Hmisc)
describe(mydata)
# n, nmiss, unique, mean, 5,10,25,50,75,90,95th percentiles
# 5 lowest and 5 highest scores
```

Using the [pastecs](#) package

```
library(pastecs)
stat.desc(mydata)
# nbr.val, nbr.null, nbr.na, min max, range, sum,
# median, mean, SE.mean, CI.mean, var, std.dev, coef.var
```



$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n-1}$$

- Vážený (z frekv. tabuľky):

$$s^2 = \frac{1}{n-1} \sum_{i=1}^m (x_i - \bar{x})^2 n_i = \frac{(x_1 - \bar{x})^2 n_1 + \dots + (x_m - \bar{x})^2 n_m}{n-1}$$

$$s = \sqrt{s^2}$$

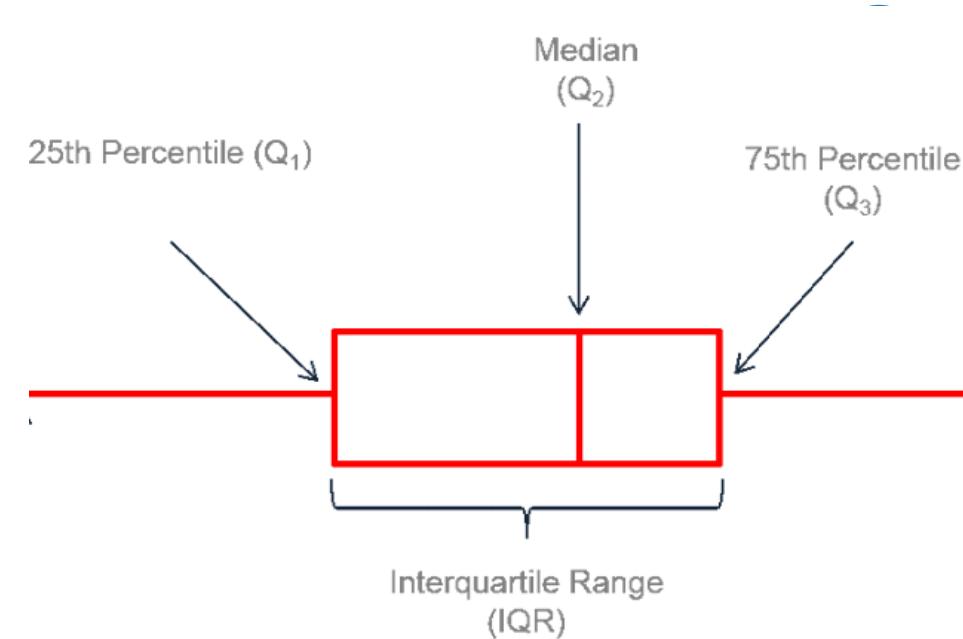
$$SE = \frac{s}{\sqrt{n}}$$

$$R = x_{\max} - x_{\min}$$

MIERY

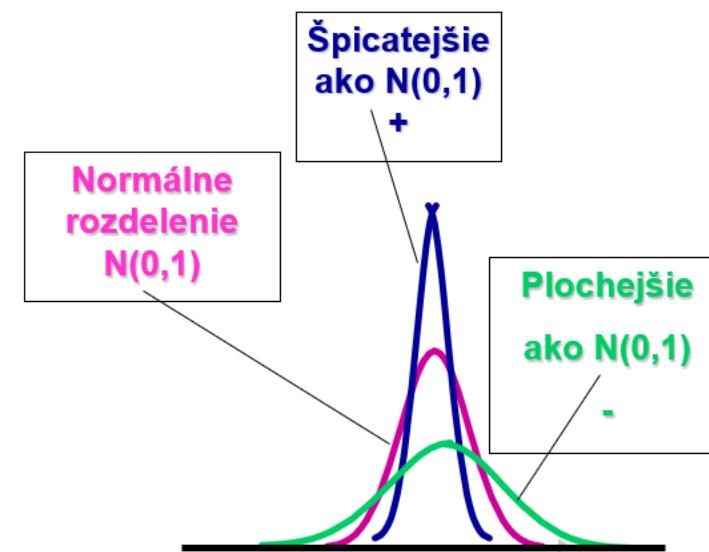
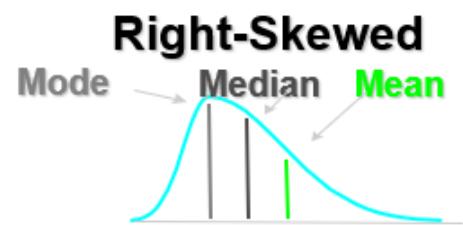
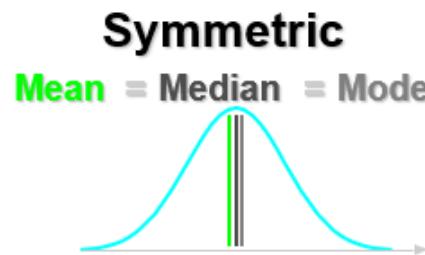
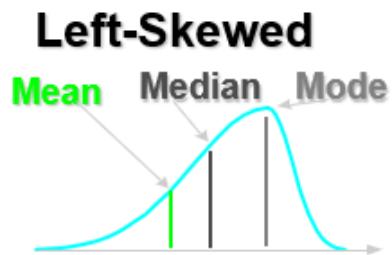
VARIABILITY

- Rozptyl
- Štandardná odchylka
- Štandardná chyba
- Variačné rozpäťie
- Medzikvartilové rozpäťie
- Variačný koeficient



$$V_k = \frac{s}{x} \cdot (100)$$

MIERY



TVARU

- Šikmost' (Skewness)
- Špicatost' (Kurtosis)

```
# calculate skewness
install.packages("moments")
library(moments)
skewness(x)

kurtosis(x)
```

Tvar (Shape) & Box Plot

Koeficient Šiknosti (Skew = Skewness), označenie γ_1

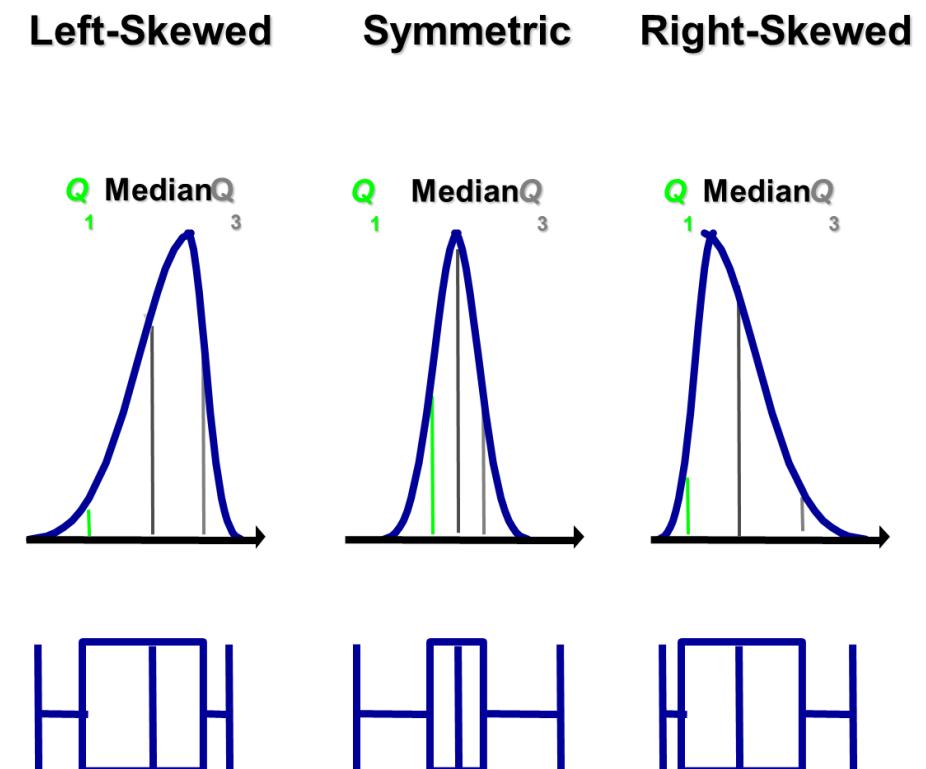
- Zošikmenie Doľava → Kladný
- Zošikmenie Doprava → Záporný
- Symetria → 0

$$\gamma_1 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{s^3}$$

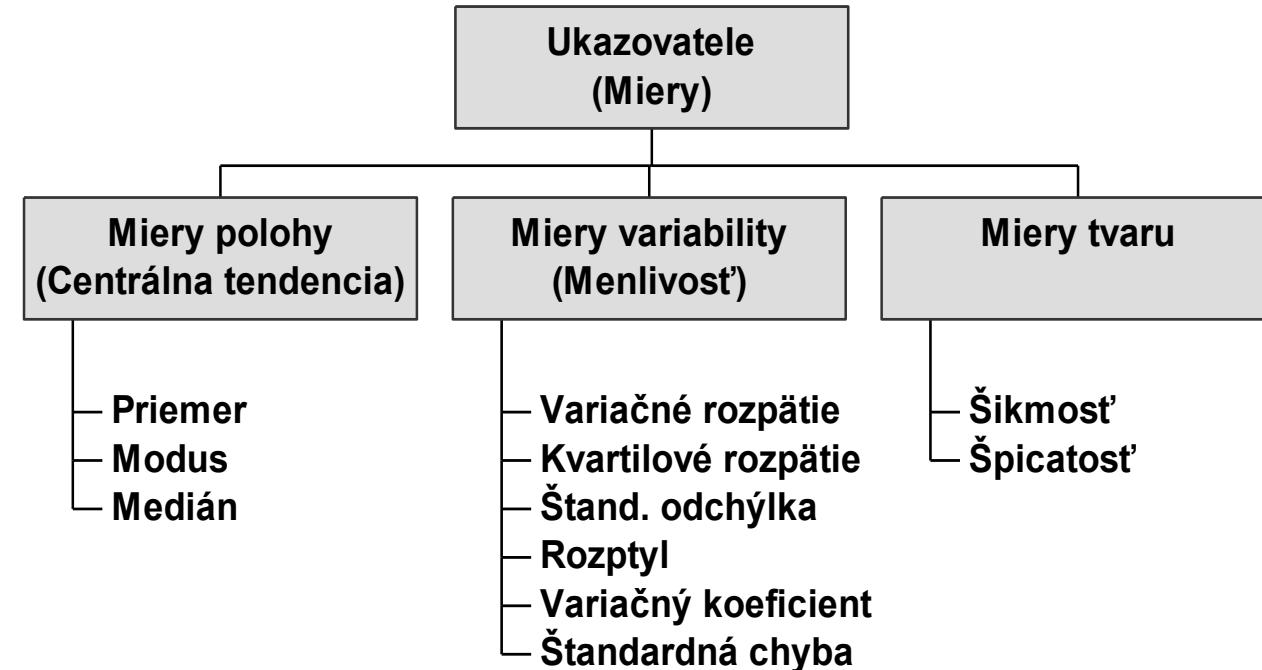
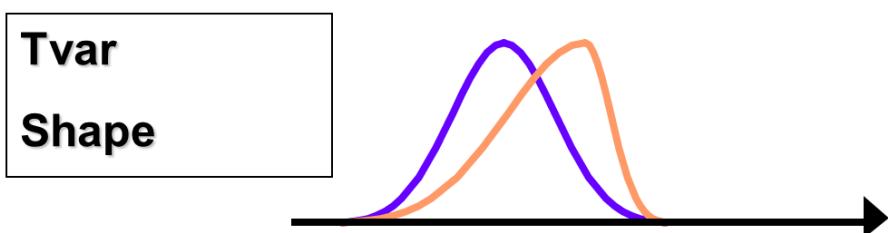
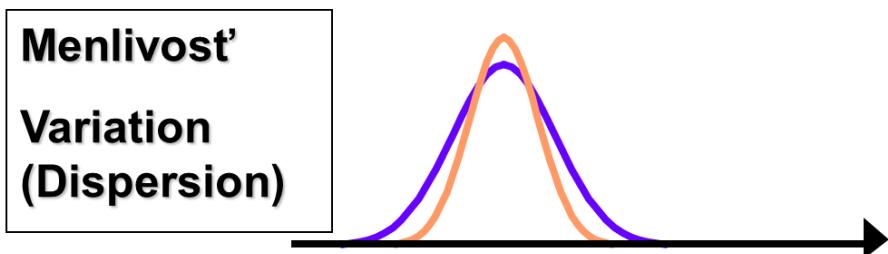
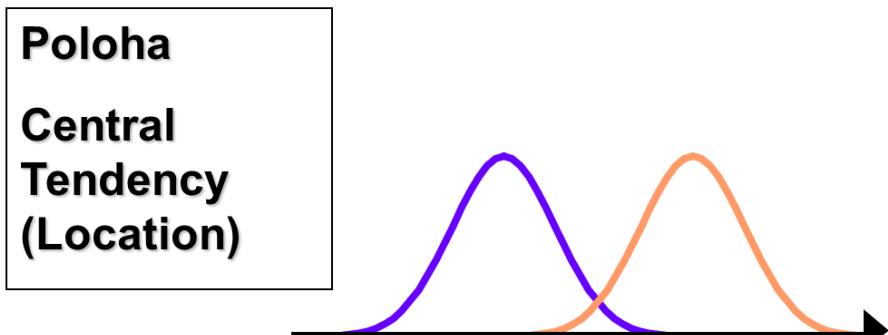
Koeficient Špicatosti (Kurt = Kurtosis) označenie γ_2

$$\gamma_2 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{s^4} - 3$$

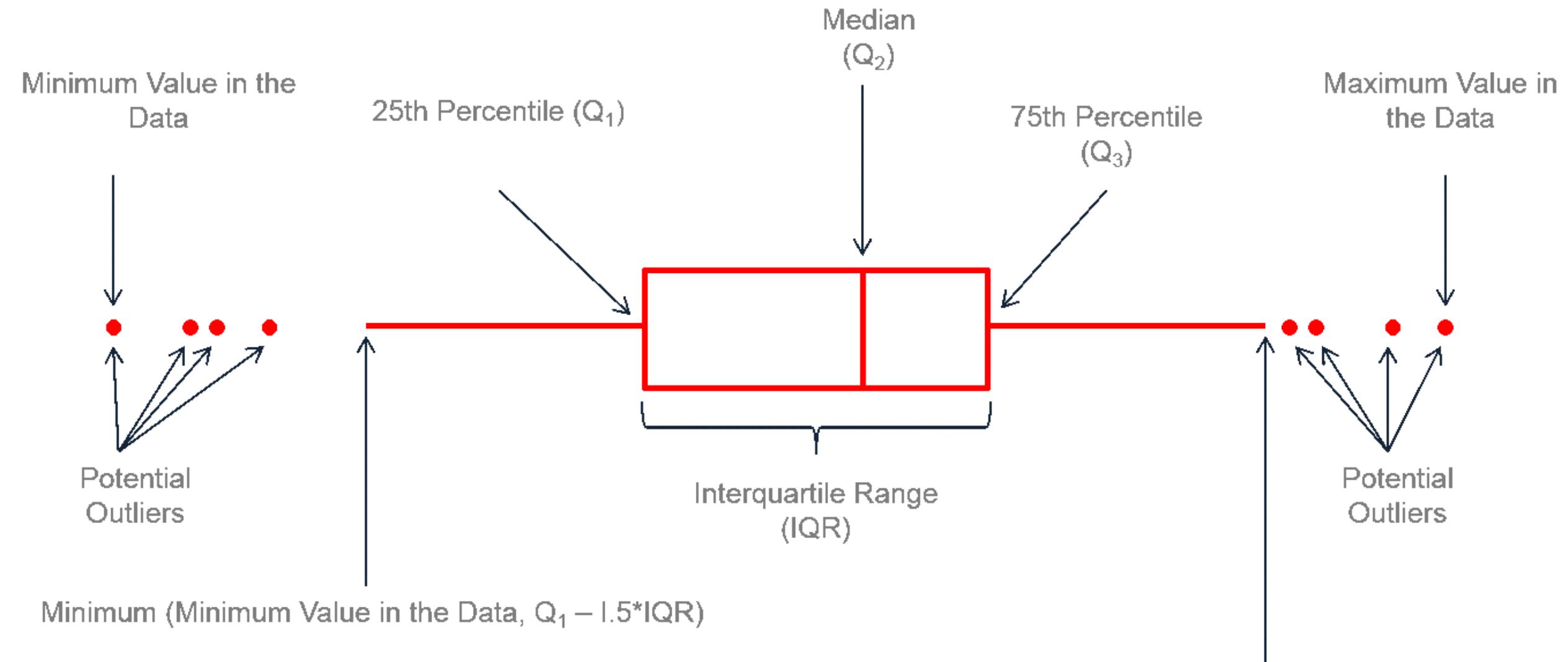
- Normálne rozdelenie → 0
- Špicatejšie → Kladný
- Plochejšie → Záporný



Ukazovatele (Miery) Vlastnosti Kvantitatívnych Dát



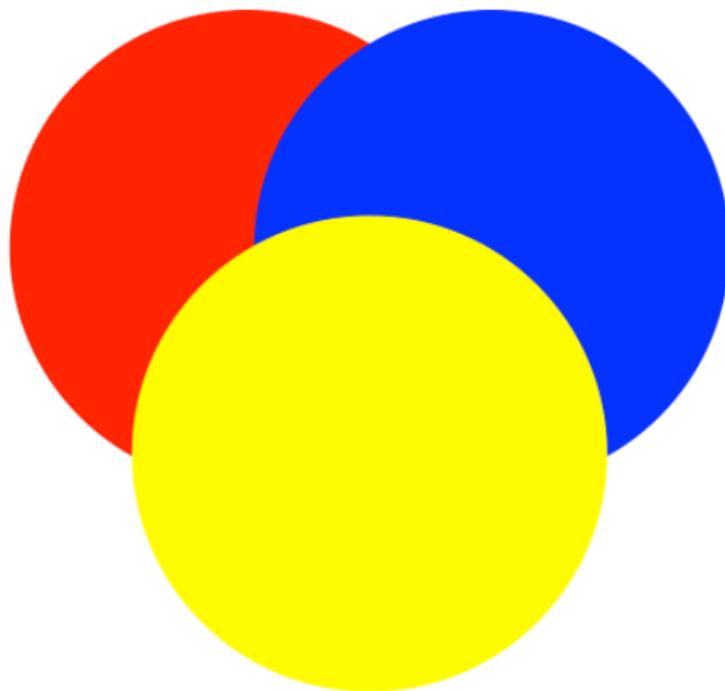
Box plot (krabicový) graf



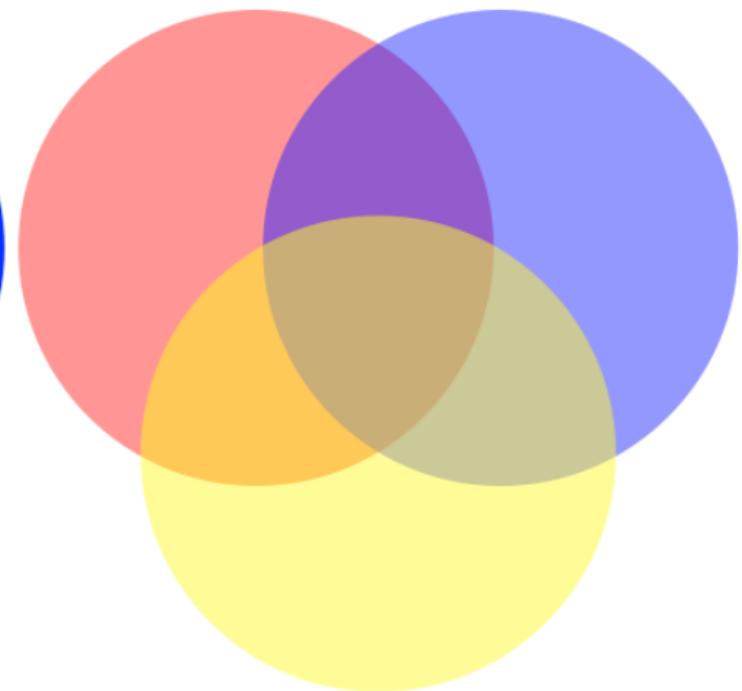
Témy a Farby v Jazyku R



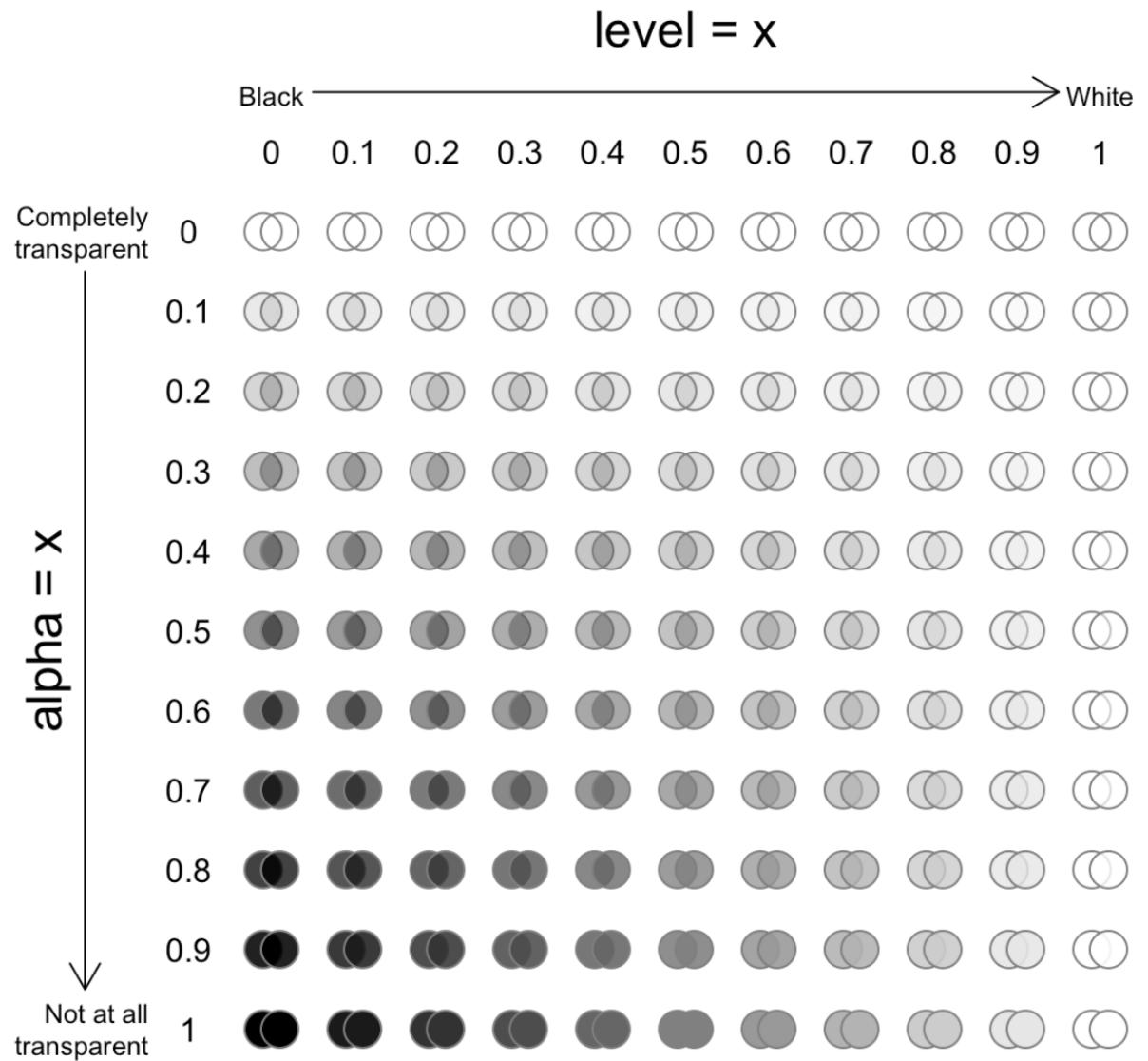
Standard



Transparent



Transparentnost' (alpha)



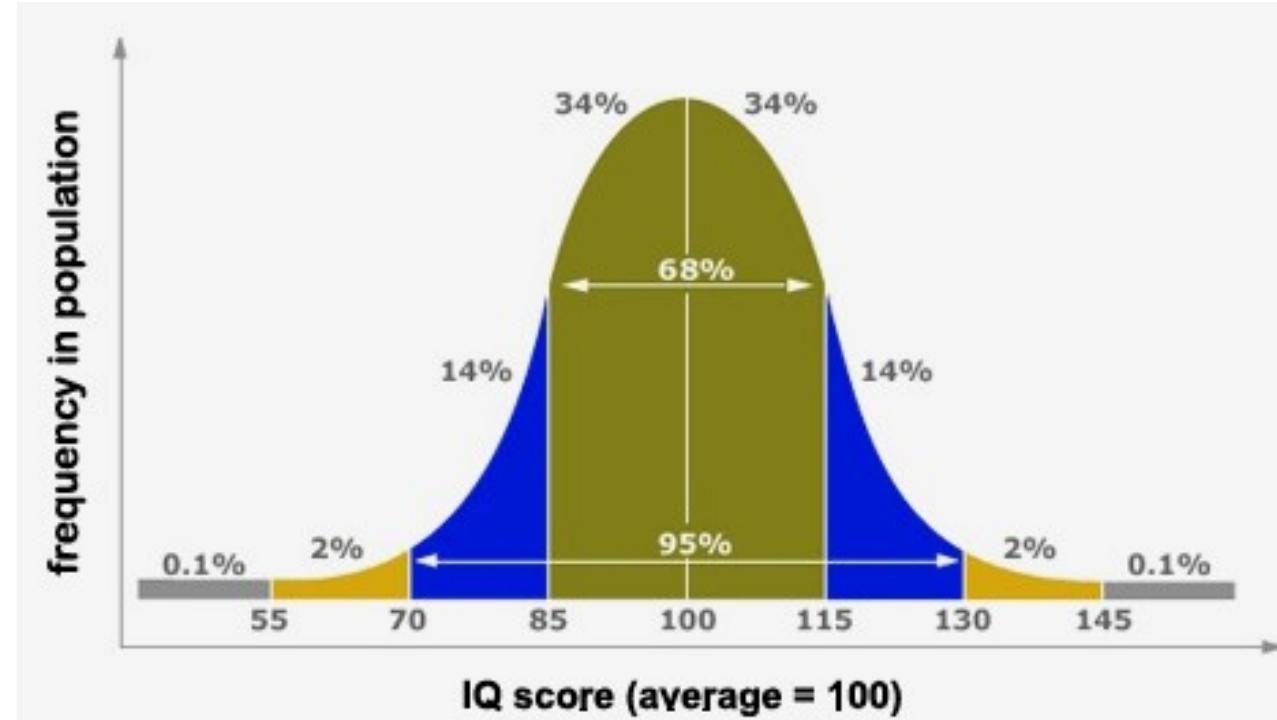


Normálne Rozdelenie (Gaussovo rozdelenie)

- Meracie prístroje (zákon chýb)
- Výrobky v továrni
- Fyzikálne a Technické Veličiny
- IQ
- Hmotnosť, výška Ľudí
- Objemy neopracovaných vzoriek
- Prirodzená Remantná Magnetizácia
- Väčšine prípadov má normálne rozdelenie, alebo rozdelenie z neho odvodené (t , F , Chi-kvadrát rozdelenie)
- Veľa premenných má normálne rozdelenie, čo sa stalo príčinou toho, že normálne rozdelenie sa považuje za všeobecnú črtu skúmanej reality
- Problém nastane, keď sa snažíme použiť takýto test na dátu, u ktorých sa normálne rozdelenie nepotvrdí

Normálne Rozdelenie

- Cieľom štatistického testovania je odhadnúť vzťah medzi premennými, teda pomer vysvetlenej varácie jednej premennej pomocou druhej voči celkovej variácii tej premennej. Teda, až do akej miery, vyjadrenej v %, sa zmena prvej premennej dá vysvetliť zmenou druhej a naopak
- Vo väčšine prípadov poznáme tvar funkcie a vieme určiť významnosť pre nás nález v príslušnej vzorke.
- Väčšina týchto funkcií súvisí s funkciou nazývanou 'normálna'
- Jej krivka, nazývaná aj **Gaussova krivka**, má tvar **zvona** a je funkciou iba 2 parametrov: **priemeru a štandardnej odchýlky**
- Normálne rozdelenie (distribúcia) reprezentuje jednu z empiricky verifikovaných práv o 'väseobecnej povahе reality'
- O normálnom rozdelení platí: **68 % populácie má hodnotu meraného znaku ležiacu v intervale priemer ± 1 krát štandardná odchýlka, 95 % v intervale priemer ± 2 krát štandardná odchýlka**
- Inými slovami, hodnota, ktorá neleží v tomto intervale má relatívnu frekvenciu 5% alebo menej



Aké IDE použiť?



Integrated Development Environment

Chceme úplne všetko!



Individual Edition is now

ANACONDA DISTRIBUTION

The world's most popular open-source Python distribution platform

Anaconda Distribution

[Download !\[\]\(850c55d3eee6d137666621e1273ca036_img.jpg\)](#)

For Windows
Python 3.9 • 64-Bit Graphical Installer • 510 MB

Get Additional Installers





Open Source

Access the open-source software you need for projects in any field, from data visualization to robotics.



User-friendly

With our intuitive platform, you can easily search and install packages and create, load, and switch between environments.

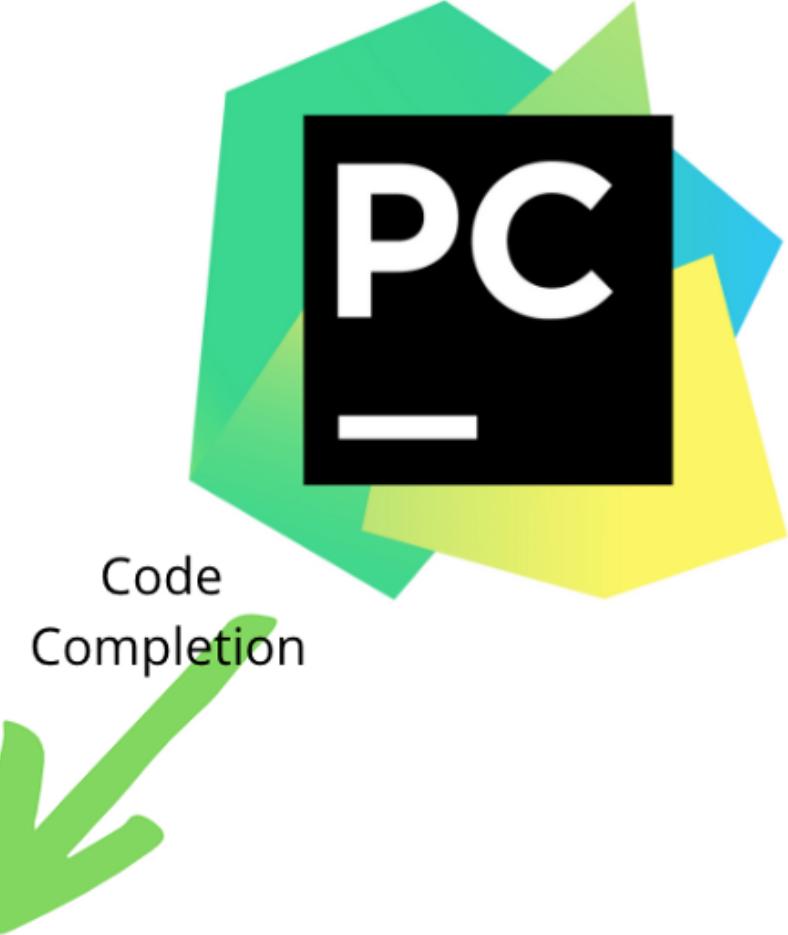
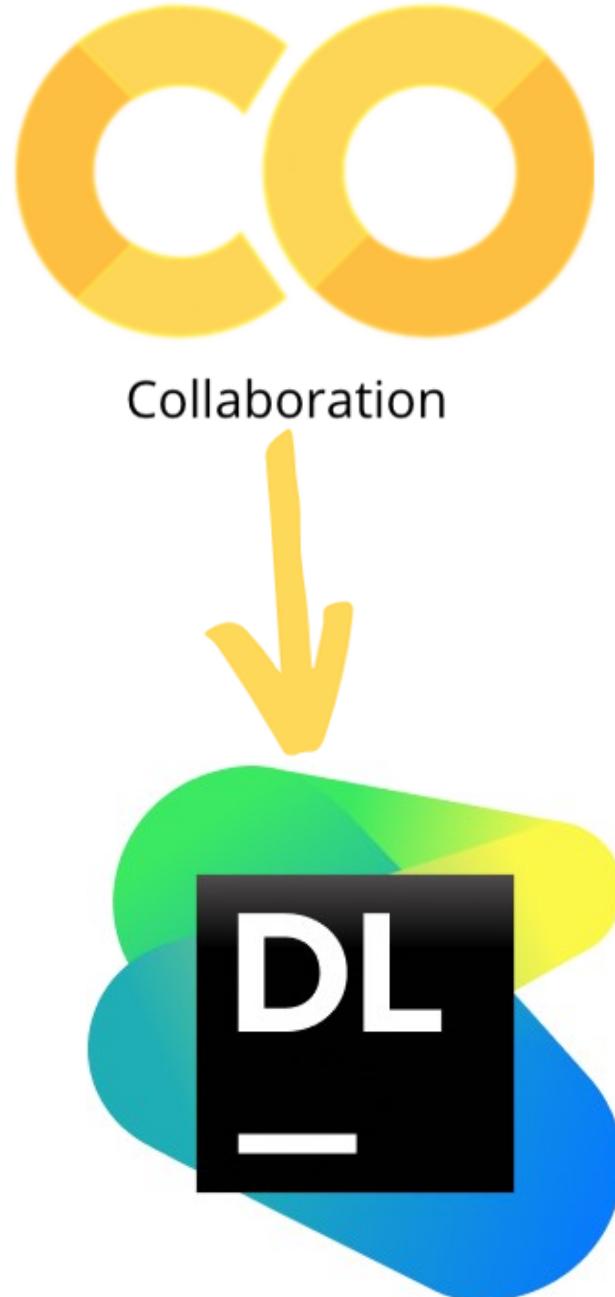


Trusted

Our securely hosted packages and artifacts are methodically tested and regularly updated.



Flexibility





MR Home ▾

[Your notebooks](#)[Recent activity](#)[Shared with you](#)[Your favorites](#)[Trash](#)[Workspace files](#)[Reports](#)[Scheduled notebooks](#)[Databases](#)[Cloud storages](#)[Environment variables](#)[Settings](#)[+ New notebook](#)

Your notebooks

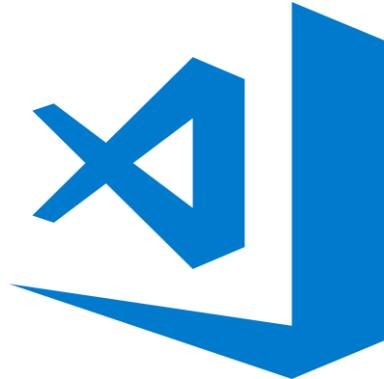
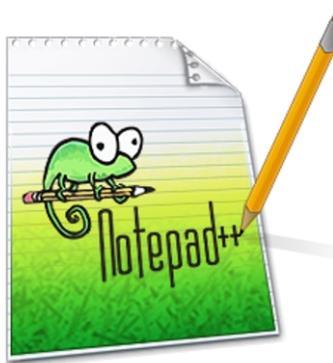
Name	Last modified	Author	
Python II. Mierne Pokrocily	26 Oct 2022 00:34	Miroslav Reiter	...
Python I. Zaciatočník	26 Oct 2022 00:34	Miroslav Reiter	...
Python 5 Automatizácia Úloh	09 May 2023 12:56	Miroslav Reiter	...
Python 4 Databázy	25 Apr 2023 23:15	Miroslav Reiter	...
Python 3 SVD	26 Oct 2022 00:32	Miroslav Reiter	...
Python 2 OOP	26 Oct 2022 00:33	Miroslav Reiter	...
Kurz SAV R Základy	10 Nov 2022 10:09	Miroslav Reiter	...
Kurz SAV R Spracovanie Udajov	13 Apr 2023 00:59	Miroslav Reiter	...
Kurz SAV Python Zaklady	15 Apr 2023 11:39	Miroslav Reiter	...
Getting started tutorial	29 Sep 2022 22:42	Miroslav Reiter	...





```
print("Sorry, we are down for maintenance")
print("We'll be back shortly")
```

Aký Editor Použiť?



```
:::  
iLE880j. :jD888880j:  
.LGitE888D.f8GjjjL8888E;  
iE :8888Et. .G888.  
;i E888, ,8888,  
D888, :8888:  
D888, :8888:  
D888, :8888:  
D888, :8888:  
888W, :8888:  
W88W, :8888:  
W88W: :8888:  
DGGD: :8888:  
:8888:  
:W888:  
:8888:  
E888i  
tW88D
```



IDE ≠ Editor

Balíčky (Inštalácia a Načítanie)

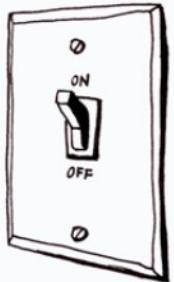
Installing a package

```
install.packages('my.package')
```

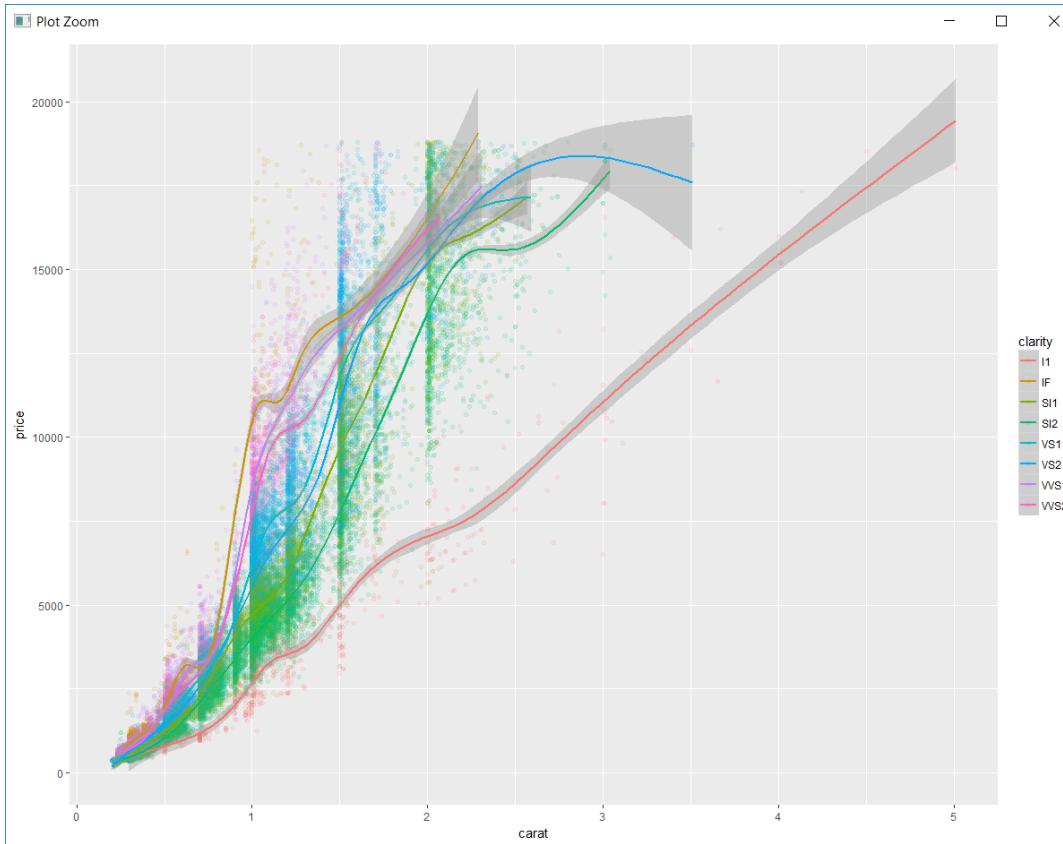


Loading a package

```
library('mypackage')
```



Inštalácia Balíčkov (packages)



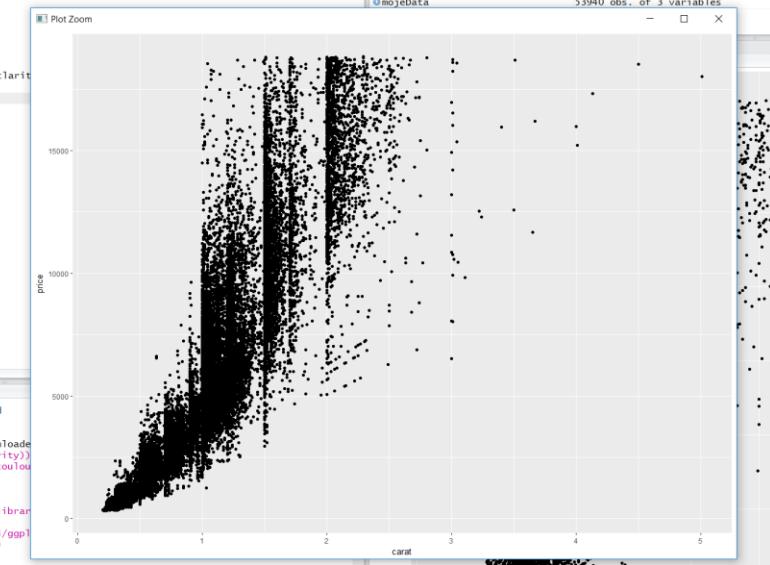
```
mojeData <- read.csv(file.choose())
install.packages("ggplot2")
ggplot(data = mojeData, aes(x=carat, y=price, colour=clarity))
geom_point()
```

The screenshot shows the RStudio interface. On the left, the code editor displays a script with the following content:

```
mojeData <- read.csv(file.choose())
install.packages("ggplot2")
ggplot(data = mojeData, aes(x=carat, y=price, colour=clarity))
geom_point()
```

On the right, the "Console" tab shows the output of the commands:

```
package 'ggplot2' successfully unpacked and MD5 sums checked
The downloaded binary packages are in
  C:/Users/Miroslav/AppData/Local/Temp/RtmpMPCJr4/downloade
> ggplot(data = mojeData, aes(x=carat, y=price, colour=clarity))
Error in ggplot(data = mojeData, aes(x = carat, y = price, coulo
could not find function "ggplot"
> #> #> #>
> install.packages("ggplot2")
Installing package into 'C:/Users/Miroslav/Documents/R/win-library
(as 'lib' is unspecified)
trying URL 'https://cran.rstudio.com/bin/windows/contrib/3.4/ggp
Content type: 'application/zip' length 2784631 bytes (2.7 MB)
downloaded 2.7 MB
```



Import Dát a Balíčkov

The screenshot shows the RStudio interface with two code files and two panes.

Hello.R:

```
1
2
3
4
5
6 mojeData <- read.csv(file.choose())
7
8 install.packages("ggplot2")
9
10 ggplot(data = mojeData, aes(x=carat, y=price, colour=clarity))
#geom_point(alpha=0.1)
11
12
13
```

Data.R:

```
1
2
3
4
5 mojeData      53940 obs. of 3 variables
```

Environment Pane:

- File: Hello.R*
- File: Data.R*
- Source on Save
- Run
- Source

Packages Pane:

- Files
- Plots
- Packages** (selected)
- Help
- Viewer

Name	Description	Version
dplyr	A Grammar of Data Manipulation	0.7.2
Formula	Extended Model Formulas	1.2-2
ggplot2	Create Elegant Data Visualisations Using the Grammar of Graphics	2.2.1
glue	Interpreted String Literals	1.1.1
gttable	Arrange 'Grobs' in Tables	0.2.0
labeling	Axis Labeling	0.3
lazyeval	Lazy (Non-Standard) Evaluation	0.2.0



dplyr

Overview

dplyr is a grammar of data manipulation, providing a consistent set of verbs that help you solve the most common data manipulation challenges:

- `mutate()` adds new variables that are functions of existing variables
- `select()` picks variables based on their names.
- `filter()` picks cases based on their values.
- `summarise()` reduces multiple values down to a single summary.
- `arrange()` changes the ordering of the rows.

These all combine naturally with `group_by()` which allows you to perform any operation “by group”. You can learn more about them in `vignette("dplyr")`. As well as these single-table verbs, dplyr also provides a variety of two-table verbs, which you can learn about in `vignette("two-table")`.

LINKS

[View on CRAN](#)

[Browse source code](#)

[Report a bug](#)

[Learn more](#)

LICENSE

[Full license](#)

[MIT + file LICENSE](#)

COMMUNITY

[Contributing.guide](#)

[Code of conduct](#)

[Getting help](#)

CITATION

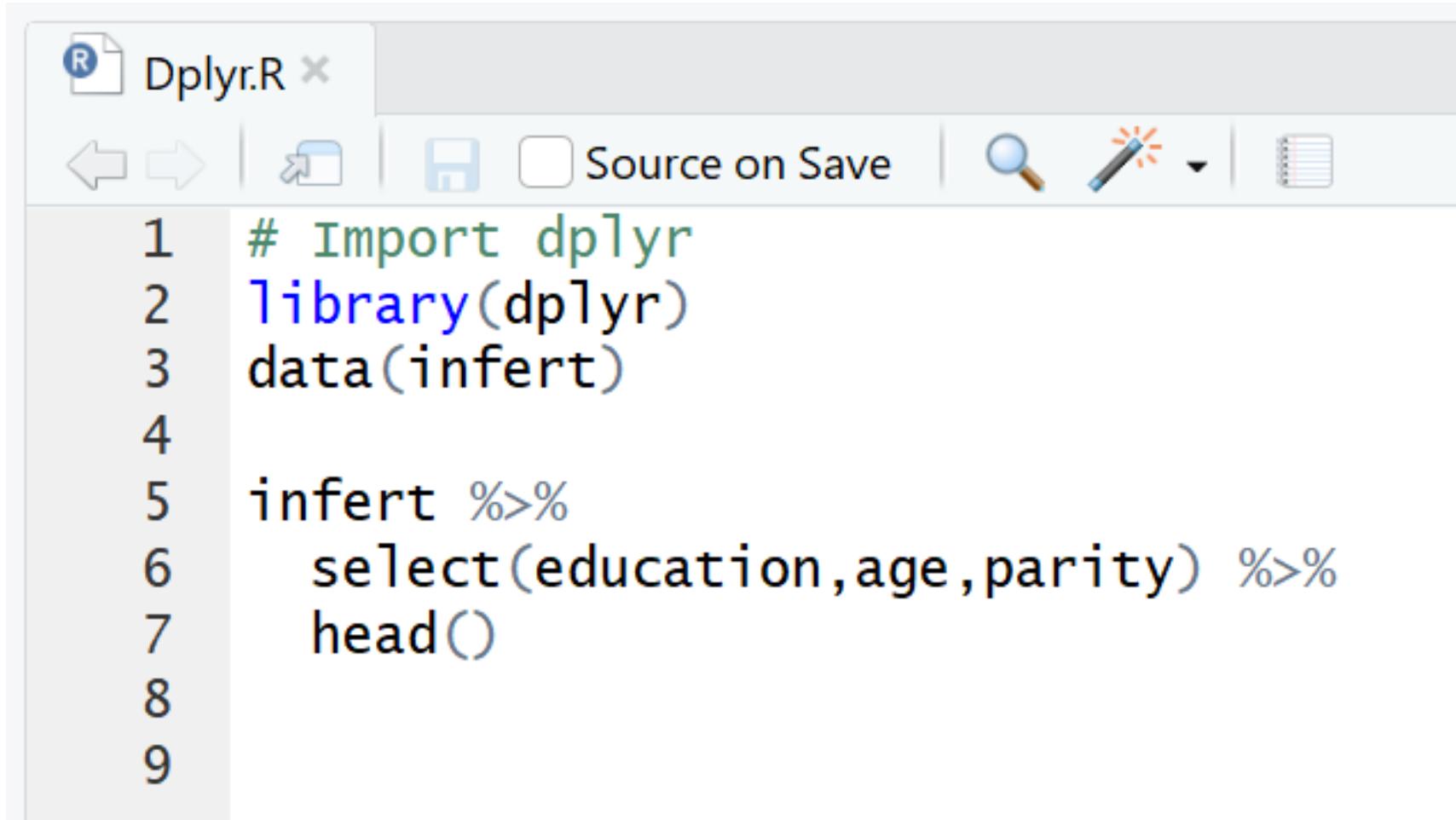
[Citing dplyr](#)

DEVELOPERS

[Hadley Wickham](#)

Author, maintainer

Výber Dát (Select)



The screenshot shows an RStudio interface with a code editor window. The window title is "Dplyr.R". The code editor contains the following R script:

```
1 # Import dplyr
2 library(dplyr)
3 data(infert)
4
5 infert %>%
6   select(education,age,parity) %>%
7   head()
```

Selecting Columns of Data

When you're dealing with large datasets, it can be confusing to look at many different columns of data at once. Here, select the `President`, `Date`, and `Approve` columns from the `approval_polls` dataset, and pipe the result to the `head()` function to display the results.

INSTRUCTIONS 100 XP

- Using the `select()` function from the `tidyverse` package, extract the `President`, `Date`, and `Approve` columns from the `approval_polls` dataset.

 Take Hint (-30 XP)

```
1 # Import dplyr
2 library(dplyr)
3
4 # Select President, Date, and Approve from approval_polls
5 approval_polls %>%
6   select(President, Date, Approve) %>%
7   head()
```

 Run Code  Submit Answer

R CONSOLE SLIDES

```
> # Import dplyr
> library(dplyr)
>
> # Select President, Date, and Approve from approval_polls
> approval_polls %>%
  select(President, Date, Approve) %>%
  head()
```

	President	Date	Approve
1	Trump	12/12/2017	36
2	Trump	12/9/2017	36
3	Trump	12/6/2017	37
4	Trump	12/3/2017	35
5	Trump	11/30/2017	34
6	Trump	11/27/2017	37

Atomické Vektory

- Commonly known as **vectors**
- **Homogeneous** data structure

Character vector	Numeric vector	Integer vector	Logical vector	Complex vector
{ "A", "c" }	{ 4.36, 7.42 }	{ 3,5 }	{ True, False }	{ 1 + 7i , 8 - 2i }

Zoznam premenných

```
> ls(all.names = T)
[1] ".Random.seed"      "a"           "A"           "B"
[5] "c"                 "jeskladom"    "komplexnecislo" "meno"
[9] "menospolocnosti"   "mzdazamestnanca" "priezvisko"   "sprava"
[13] "vekklienta"        "vysledok"     "x2"          "x3"
> ls(all.names = T,sorted = T)
[1] ".Random.seed"      "a"           "A"           "B"
[5] "c"                 "jeskladom"    "komplexnecislo" "meno"
[9] "menospolocnosti"   "mzdazamestnanca" "priezvisko"   "sprava"
[13] "vekklienta"        "vysledok"     "x2"          "x3"
> ls(all.names = T,sorted = F)
[1] "sprava"            "jeskladom"    "A"           "B"
[5] "c"                 "vysledok"     "a"           "komplexnecislo"
[9] "meno"               ".Random.seed"  "mzdazamestnanca" "priezvisko"
[13] "vekklienta"        "menospolocnosti" "x2"          "x3"
> ls(all.names = T,sorted = T, pattern = "a")
[1] ".Random.seed"      "a"           "jeskladom"    "mzdazamestnanca"
[5] "sprava"            "vekklienta"
> ls(all.names = T,sorted = T, pattern = "m*")
[1] ".Random.seed"      "a"           "A"           "B"
[5] "c"                 "jeskladom"    "komplexnecislo" "meno"
[9] "menospolocnosti"   "mzdazamestnanca" "priezvisko"   "sprava"
[13] "vekklienta"        "vysledok"     "x2"          "x3"
```

Príkaz ls() | objects()

Zoznam Objektov a Funkcií

```
> ls(all.names = T)
[1] ".Random.seed"      "a"           "A"           "B"
[5] "C"                 "jeskladom"    "komplexneCislo" "meno"
[9] "menospolocnosti"  "mzdazamestnanca" "priezvisko"   "sprava"
[13] "vekklienta"        "vysledok"     "x2"          "x3"
> ls(all.names = T,sorted = T)
[1] ".Random.seed"      "a"           "A"           "B"
[5] "C"                 "jeskladom"    "komplexneCislo" "meno"
[9] "menospolocnosti"  "mzdazamestnanca" "priezvisko"   "sprava"
[13] "vekklienta"        "vysledok"     "x2"          "x3"
> ls(all.names = T,sorted = F)
[1] "sprava"            "jeskladom"    "A"           "B"
[5] "C"                 "vysledok"     "a"           "komplexneCislo"
[9] "meno"               ".Random.seed"  "mzdazamestnanca" "priezvisko"
[13] "vekklienta"        "menospolocnosti" "x2"          "x3"
> ls(all.names = T,sorted = T, pattern = "a")
[1] ".Random.seed"      "a"           "jeskladom"    "mzdazamestnanca"
[5] "sprava"            "vekklienta"
> ls(all.names = T,sorted = T, pattern = "m*")
[1] ".Random.seed"      "a"           "A"           "B"
[5] "C"                 "jeskladom"    "komplexneCislo" "meno"
[9] "menospolocnosti"  "mzdazamestnanca" "priezvisko"   "sprava"
[13] "vekklienta"        "vysledok"     "x2"          "x3"
```

Príkaz `builtins()` | `ls(baseenv(), all = TRUE)`

Dátové Štruktúry

Homogénne typy

- Rovnaké typy položiek
 - Atomické Vektory (Vectors)
 - Matice (Matrices)
 - Polia (Arrays)

Heterogénne typy

- Rôzne typy položiek
 - Zoznamy (Lists)
 - Dátové Rámce (DataFrames)

Ďalšie Dátové Typy (Dátové Štruktúry)

Matice alebo Polia

- Sú viacozmerné zovšeobecnenia vektorov
- V skutočnosti sú to vektory, ktoré môžu byť indexované dvomi alebo viacerými indexmi

Faktory

- Kompaktné spôsoby spracovania kategorických údajov

Zoznamy

- Všeobecná forma vektora, v ktorej rôzne prvky nemusia byť rovnakého typu a sú často samotné vektory alebo zoznamy

Dátové Rámce

- Sú maticové štruktúry, v ktorých môžu byť stĺpce rôznymi typmi
- Premýšľajte o dátových rámcoch ako o "dátových maticiach,, tabuľkách

Funkcie

- Sú samotné objekty v R, ktoré môžu byť uložené v pracovnom priestore projektu

Konverzie dát

	to one long vector	to matrix	to data frame
from vector	<code>c(x,y)</code>	<code>cbind(x,y)</code> <code>rbind(x,y)</code>	<code>data.frame(x,y)</code>
from matrix	<code>as.vector(mymatrix)</code>		<code>as.data.frame(mymatrix)</code>
from data frame		<code>as.matrix(myframe)</code>	

Dátumy (Dates)

```
# use as.Date( ) to convert strings to dates  
mydates <- as.Date(c("2007-06-22", "2004-02-13"))  
# number of days between 6/22/07 and 2/13/04  
days <- mydates[1] - mydates[2]
```

```
# print today's date  
today <- Sys.Date()  
format(today, format="%B %d %Y")  
"June 20 2007"
```

Formátovanie Dátumy

Symbol	Meaning	Example
%d	day as a number (0-31)	01-31
%a	abbreviated weekday	Mon
%A	unabbreviated weekday	Monday
%m	month (00-12)	00-12
%b	abbreviated month	Jan
%B	unabbreviated month	January
%y	2-digit year	07
%Y	4-digit year	2007

Konverzie Dátumy

```
# convert date info in format 'mm/dd/yyyy'  
strDates <- c("01/05/1965", "08/16/1975")  
dates <- as.Date(strDates, "%m/%d/%Y")
```

```
mydates <- as.Date(c("2007-06-22", "2004-02-13"))
```

```
# convert dates to character data  
strDates <- as.character(dates)
```

default format yyyy-mm-dd

Generovanie dát

Sekvencie

- `1:10`,
- `seq(10)`,
- `seq(1, 10, 2)`

Replikácie/opakovanie

- `rep(1:4, 2)`,
- `rep(1:4, each = 2)`

Permutácie a vzorové dáta

- `sample(1:20, 9)`

`as.Date("1998-12-17")`

Matematické funkcie

Aritmetické funkcie	Goniometrické funkcie	Hyperbolické funkcie	Množinové operácie
<ul style="list-style-type: none">• <code>log(x)</code>• <code>logb()</code>• <code>log10()</code>• <code>log2()</code>• <code>exp()</code>• <code>expm1()</code>• <code>log1p()</code>• <code>sqrt()</code>	<ul style="list-style-type: none">• <code>cos()</code>• <code>sin()</code>• <code>tan()</code>• <code>acos()</code>• <code>asin()</code>• <code>atan()</code>• <code>atan2()</code>	<ul style="list-style-type: none">• <code>cosh()</code>• <code>sinh()</code>• <code>tanh()</code>• <code>acosh()</code>• <code>asinh()</code>• <code>atanh()</code>	<ul style="list-style-type: none">• <code>union()</code>• <code>intersect()</code>• <code>setdiff()</code>• <code>setequal()</code>

Frekvenčné a Krížové Tabuľky

```
# 2-Way Frequency Table
attach(mydata)
mytable <- table(A, B) # A will be rows, B will be columns
mytable # print table

margin.table(mytable, 1) # A frequencies (summed over B)
margin.table(mytable, 2) # B frequencies (summed over A)

prop.table(mytable) # cell percentages
prop.table(mytable, 1) # row percentages
prop.table(mytable, 2) # column percentages

# 3-Way Frequency Table
mytable <- table(A, B, C)
ftable(mytable)
```

Frekvenčné a Krížové Tabuľky 3-way

```
# 3-Way Frequency Table  
mytable <- table(A, B, C)  
ftable(mytable)
```

```
# 3-Way Frequency Table  
mytable <- xtabs(~A+B+c, data=mydata)  
ftable(mytable) # print table  
summary(mytable) # chi-square test of indepedence
```

Frekvenčné a Krížové Tabuľky

```
# 2-Way Cross Tabulation  
library(gmodels)  
CrossTable(mydata$myrowvar, mydata$mycolvar)
```

Testy Nezávislosti

Test Chi-Square

- Pre dvojcestné tabuľky môžete použiť test chisq.test (mytable) na overenie nezávislosti premennej riadka a stĺpca
- Predvolene je hodnota p vypočítaná z asymptotickej distribúcie štatistických údajov testu
- Prípadne môže byť hodnota p odvodená pomocou simultánnej funkcie Monte Carlo

Fisher Exact Test

- fisher.test (x) poskytuje presný test nezávislosti
- x je dvojrozmerná kontingenčná tabuľka v maticovej forme

```
> chisq.test(zlocinci)
```

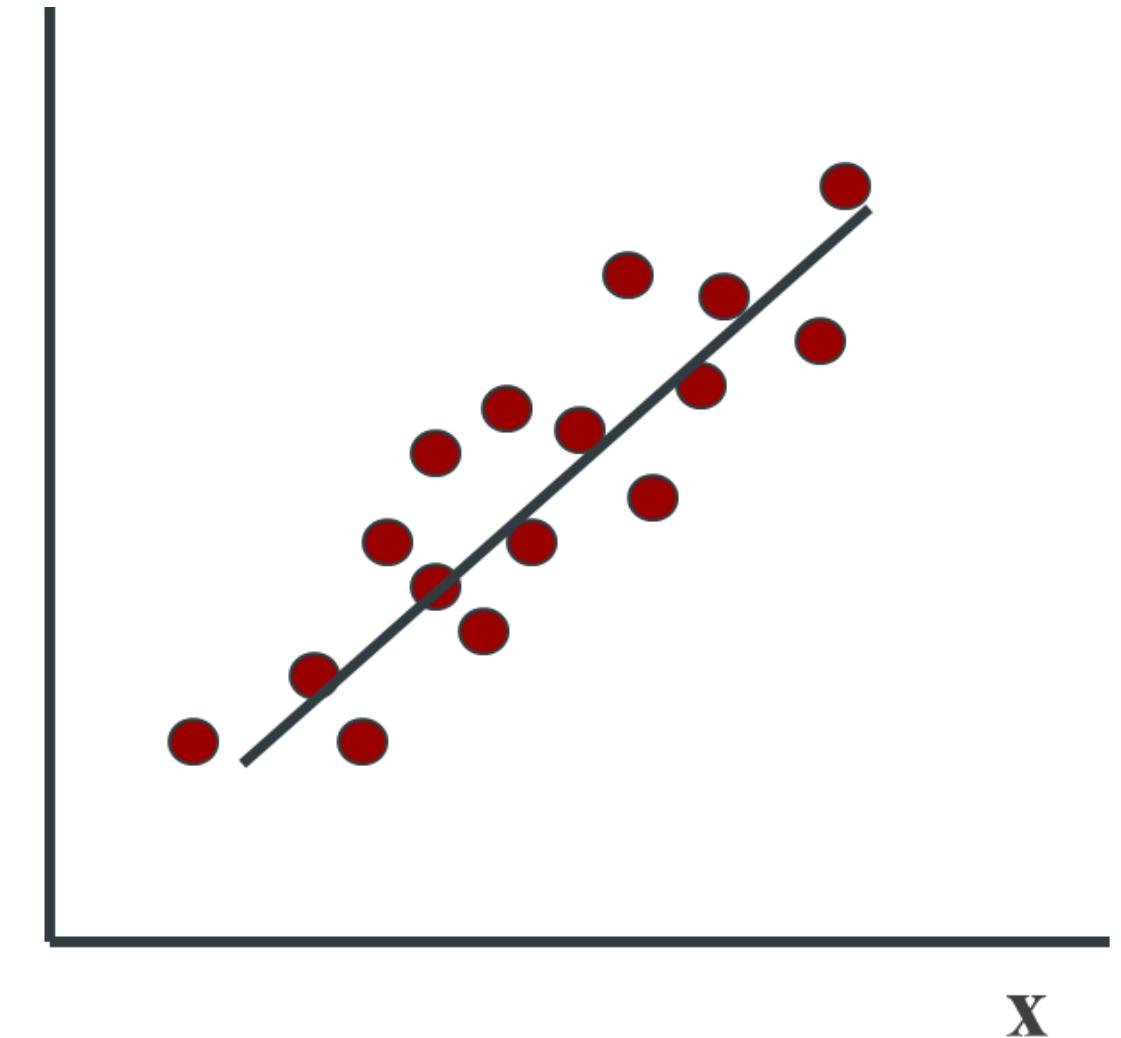
Pearson's chi-squared test

```
data: zlocinci  
X-squared = 714.85, df = 147, p-value < 2.2e-16
```

```
Warning message:  
In chisq.test(zlocinci) : chi-squared approximation may be incorrect  
> |
```

Korelačná Analýza

- Overenie vypovedacej schopnosti kvantifikovaných regresných modelov ako celku, aj jeho častí
- Výpočet **číselných charakteristik**, ktoré v koncentrovanej forme popisujú **kvalitu vypočítaných modelov**
- Požadujeme od nich, aby sa pohybovali v pevne ohraničenom intervale
- V rámci intervalu rástli s vyššou silou závislosti
- Porovnanie 2 prípadov závislosti



Korelačná Analýza

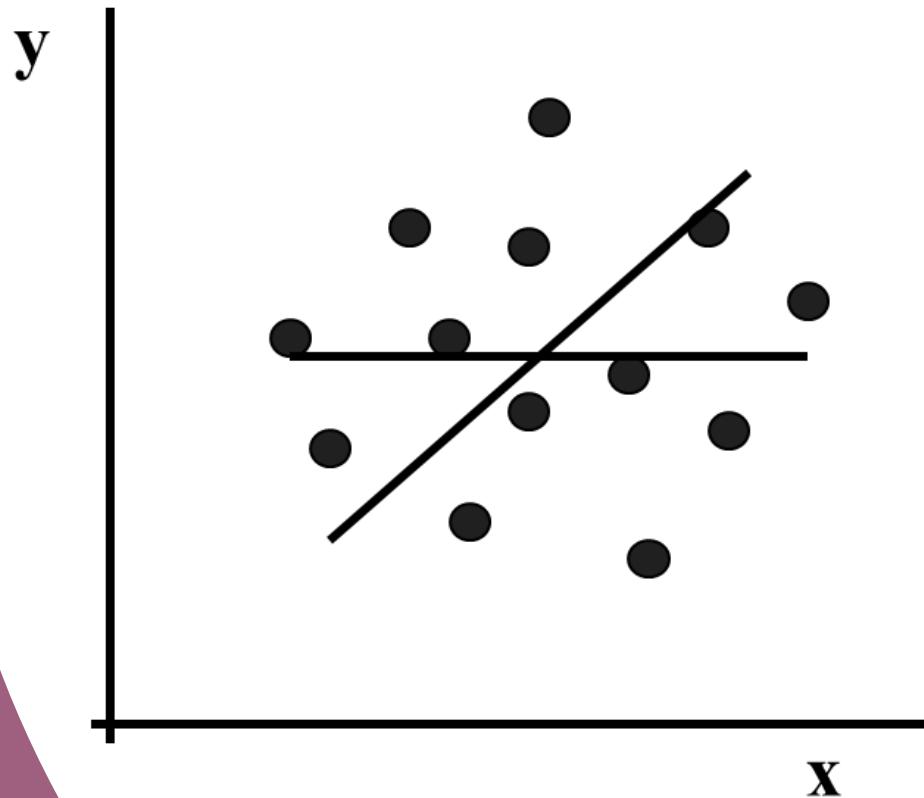
- Po úvodnom **grafickom preskúmaní** nastupuje **fáza hľadania presných štatistik**, ktoré potvrdia odhady z grafov
- Korelačné štatistiky zistujú **či medzi premennými existuje korelacia aká je sila korelácie**
- **Koreláciou nazývame vzájomný lineárny vztah – závislosť 2 premenných**
- Tento vztah môže byť:
 1. **Priamy** – s rastúcimi hodnotami 1. premennej rastú aj hodnoty 2. premennej
 2. **Nepriamy** – s rastúcimi hodnotami 1. premennej klesajú hodnoty 2. premennej
 3. Ak medzi hodnotami 2 premenných neexistuje ani priama ani nepriama lineárna závislosť, hovoríme, že sú **nekorelované**.



Korelačná Analýza

Miery tesnosti štatistickej závislosti:

- **Kovariancia** – cov_{yx}
 - len pre lineárnu závislosť
- **Koeficient Korelácie** r_{yx}
 - len pre lineárnu závislosť
- **Koeficient Determinácie** r_{yx}^2
 - len pre lineárnu závislosť
 - index korelácie i_{yx}
 - index determinácie i_{yx}^2



Kovariancia

- Miera, ktorá sa používa, aby sme potvrdili alebo vyvrátili existenciu lineárnej závislosti (korelácie)
- Zo spôsobu výpočtu možno odvodiť, kedy potvrdzuje existenciu pozitívnej, negatívnej korelácie a kedy nekorelovanosti
- Ak kovariancia potvrdí neexistenciu lineárneho vzťahu, medzi premennými môže existovať nelineárny vzťah
- Ak kovariancia potvrdí existenciu lineárneho vzťahu, môžeme merat' jeho intenzitu

$$\text{cov } xy = \frac{1}{n} \sum \left(x_i - \bar{x} \right) \left(y_i - \bar{y} \right)$$

$\text{cov } xy = 0$, medzi premennými nie je lineárny vzťah
 $\text{cov } xy > 0$, medzi premennými je priamy lin. vzťah
 $\text{cov } xy < 0$, medzi premennými je nepriamy lin. vzťah

Korelačná Analýza

- Koeficient korelácie - r_{yx}

$$r_{yx} = \frac{\text{cov}_{yx}}{s_x s_y}$$

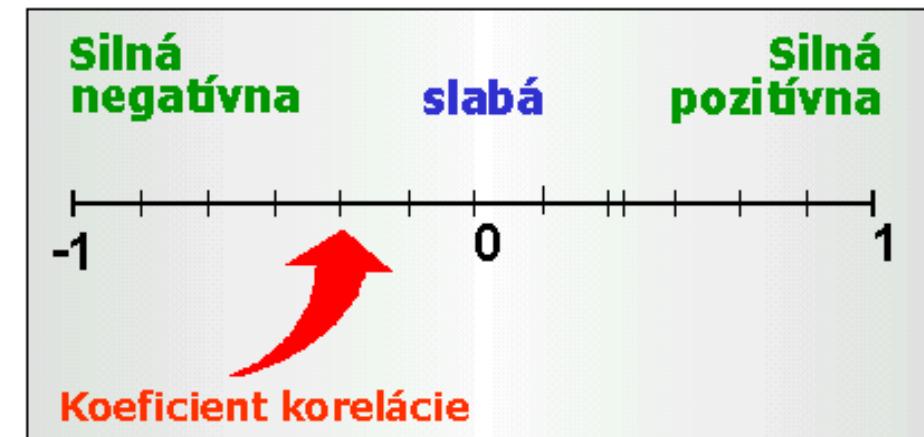
■ Hodnoty sa pohybujú v intervale: -1, 1

- ❖ $r_{yx} = -1$ – silná negatívna závislosť
- ❖ $r_{yx} = 0$ – bez závislosti
- ❖ $r_{yx} = 1$ – silná pozitívna závislosť

■ Koeficient determinácie r_{yx}^2

■ Hodnoty sa pohybujú v intervale: 0, 1

■ Udáva % vysvetlenej variability závisle premennej



Index Korelácie a Index Determinácie

- V ZSI I_{yx} odhadom z výberových údajov je i_{yx} est $I_{yx} = i_{yx}$.
- Princíp spočíva v rozklade variability závisle premennej Y

$$\sum_{j=1}^n (y_j - \bar{y})^2 = \sum_{j=1}^n (y_j' - \bar{y})^2 + \sum_{j=1}^n (y_j - y_j')^2$$

Celková
variabilita
závisle
pременнеј

Variabilita závisle
pременнеј vysvetlená
regresnou funkciou

Variabilita
nevysvetlená
regresnou funkciou –
reziuálna variabilita

Index korelácie a index determinácie

- Index korelácie i_{yx}

$$i_{yx} = \sqrt{\frac{\sum_{j=1}^n (y_j' - \bar{y})^2}{\sum_{j=1}^n (y_j - \bar{y})^2}} = \sqrt{\frac{V}{C}}$$

- Index determinácie i_{yx}^2

$$i_{yx}^2 = \frac{C - N}{C} = 1 - \frac{N}{C} = 1 - \frac{\sum_{j=1}^n (y_j - y_j')^2}{\sum_{j=1}^n (y_j - \bar{y})^2}$$

Index korelácie a Index determinácie

Index Korelácie

- Hodnoty sa pohybujú v intervale od (0,1)
- Čím sa hodnota indexu blíži k 1, tým je tesnosť závislosti vyššia a opačne

Index Determinácie

- Nadobúda hodnoty z intervalu 0 až 1
- Čím viac sa hodnota indexu blíži k 1, tým väčšia časť celkovej variability je modelom vysvetlená a naopak
- Ak sa index determinácie blíži k 0, tým menšia časť celkovej variability je vysvetlená modelom

Index determinácie

- Kritérium pri rozhodovaní o voľbe konkrétneho tvaru regresnej funkcie
- Volíme ten model, ktorý má vyšší koeficient determinácie (vyššie % vysvetlenej variability)
- Ak však majú regresné funkcie rôzny počet parametrov, je potrebné upraviť index determinácie do korigovanej podoby v tvare →
- Výrazný rozdiel medzi i^2 a $i^2_{adj.}$ indikuje, že do modelu bolo zahrnutých príliš veľa premenných

$$i^2_{\text{korig}} = 1 - \frac{(n-1) \sum_{j=1}^n (y_j - y'_j)^2}{(n-p) \sum_{j=1}^n (y_j - \bar{y})^2}$$

Korelácia

Option	Description
x	Matrix or data frame
use	Specifies the handling of missing data. Options are all.obs (assumes no missing data - missing data will produce an error), complete.obs (listwise deletion), and pairwise.complete.obs (pairwise deletion)
method	Specifies the type of correlation. Options are pearson , spearman or kendall .

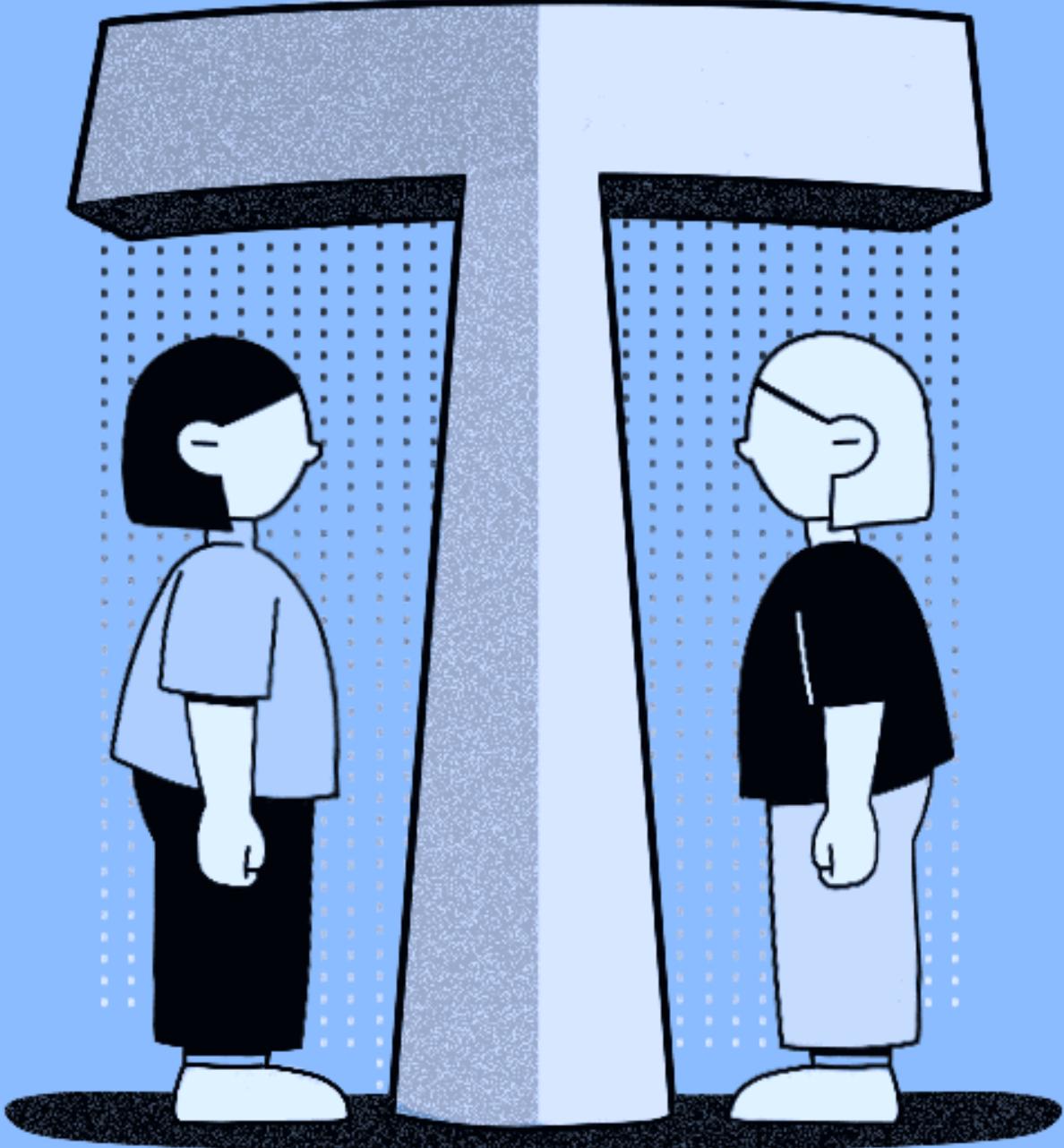
```
# Correlations/covariances among numeric variables in  
# data frame mtcars. Use listwise deletion of missing data.  
cor(mtcars, use="complete.obs", method="kendall")  
cov(mtcars, use="complete.obs")
```

Korelácia

```
# Correlations with significance levels
library(Hmisc)
rcorr(x, type="pearson") # type can be pearson or spearman

#mtcars is a data frame
rcorr(as.matrix(mtcars))

# Correlation matrix from mtcars
# with mpg, cyl, and disp as rows
# and hp, drat, and wt as columns
x <- mtcars[1:3]
y <- mtcars[4:6]
cor(x, y)
```



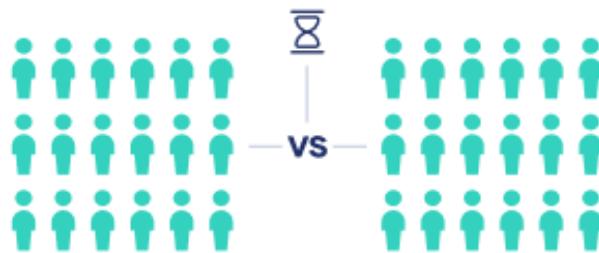
T-Test

[tē-'test]

A statistical test used to compare the means of two groups of data.

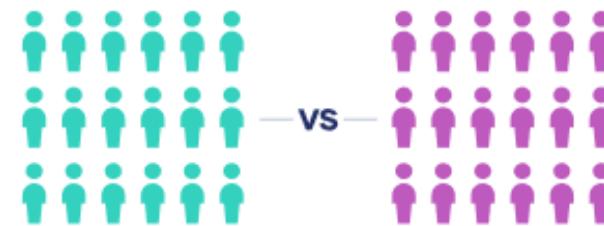
Typy t-testov

Paired-samples t test



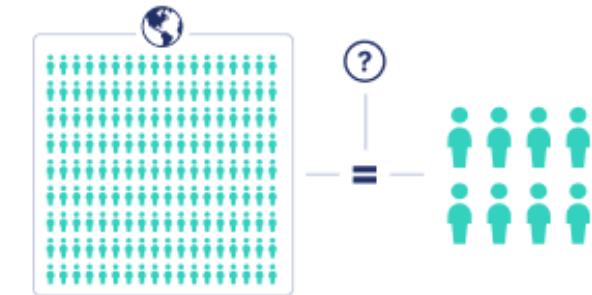
Investigate whether there's a difference within a group between two points in time (within-subjects).

Independent-samples t test



Investigate whether there's a difference between two groups (between-subjects).

One-sample t test

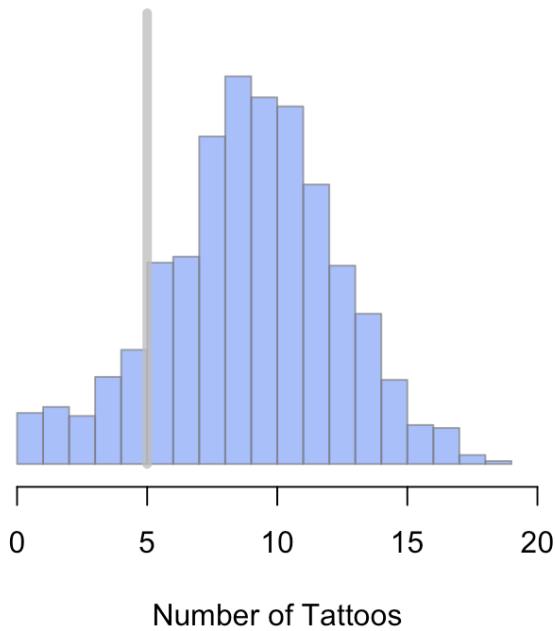


Investigate whether there's a difference between a group and a standard value or whether a subgroup belongs to a population.

Typy t-testov

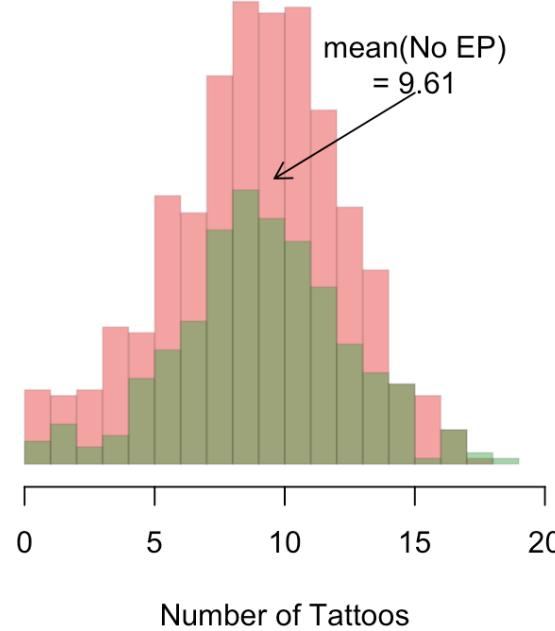
1-Sample t-test

Null Hypothesis
Mean = 5

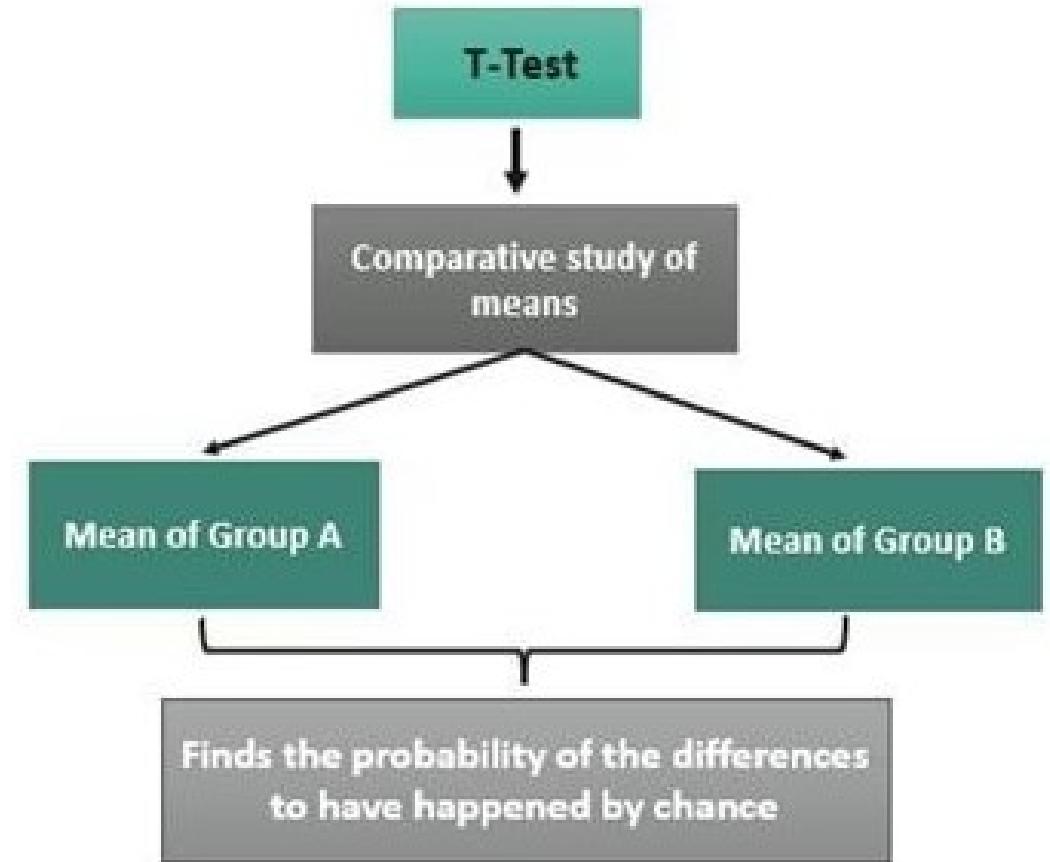


2-Sample t-test

mean(EP) = 9.34
mean(No EP) = 9.61



T-Test



T-testy

```
# independent 2-group t-test
t.test(y~x) # where y is numeric and x is a binary factor
```

$$t = \frac{m_A - m_B}{\sqrt{s^2} + \sqrt{s^2}} \sqrt{\frac{n_A}{n_B}}$$

```
# independent 2-group t-test
t.test(y1,y2) # where y1 and y2 are numeric
```

```
# paired t-test
t.test(y1,y2,paired=TRUE) # where y1 & y2 are numeric
```

$$t = \frac{m}{s/\sqrt{n}}$$

```
# one sample t-test
t.test(y,mu=3) # Ho: mu=3
```

$$t = \frac{m - \mu}{s/\sqrt{n}}$$



var.equal = TRUE



Old Shatterhand a T-Test (1-výberový)



T-test 2 sample

- t.test(Petal.Length ~ Species, data = flower.data)
- Hodnota t: -33,719. Všimnite si, že je negatívny; toto je fajn! Vo väčšine prípadov nás zaujíma len absolútна hodnota rozdielu alebo vzdialenosť od 0. Nezáleží na tom, ktorým smerom.
- Stupeň voľnosti: 30,196. Stupeň voľnosti súvisia s veľkosťou vašej vzorky a ukazujú, koľko „bezplatných“ údajových bodov je dostupných vo vašom teste na porovnávanie. Čím väčšie sú stupne voľnosti, tým lepšie bude fungovať váš štatistický test.
- Hodnota p: 2,2e-16 (t.j. 2,2 s 15 nulami). Toto popisuje pravdepodobnosť, že by ste náhodou videli takú veľkú hodnotu t, ako je táto.
- Vyjadrenie alternatívnej hypotézy (Ha). V tomto teste Ha znamená, že rozdiel nie je 0.
- 95 % interval spoľahlivosti. Toto je rozsah čísel, v ktorom bude skutočný rozdiel v priemeroch 95 % času.
- Toto sa dá zmeniť z 95 %, ak chcete väčší alebo menší interval, ale veľmi bežne sa používa 95 %.
- Priemerná dĺžka okvetných lístkov pre každú skupinu.

Welch Two Sample t-test

```
data: Petal.Length by Species
t = -33.719, df = 30.196, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-4.331287 -3.836713
sample estimates:
mean in group setosa mean in group virginica
1.456           5.540
```

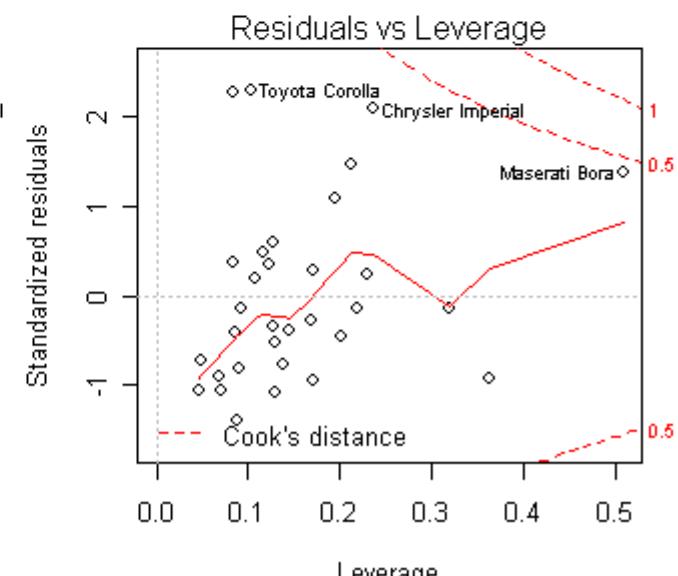
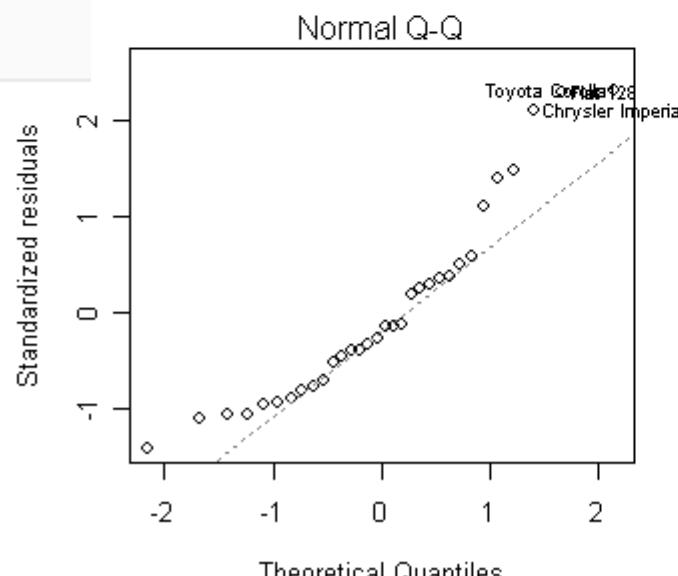
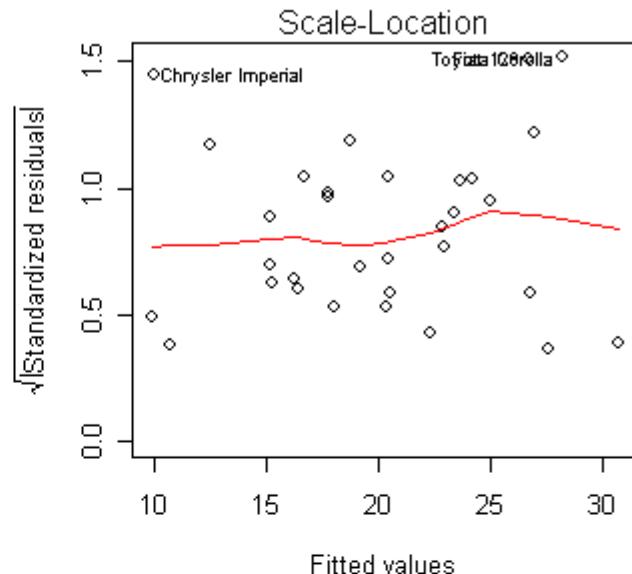
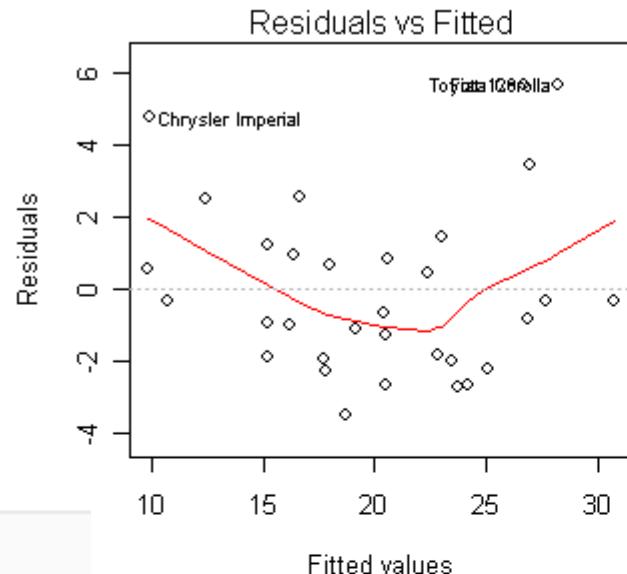
Viacnásobná (Lineárna) Regresia

```
# Multiple Linear Regression Example  
fit <- lm(y ~ x1 + x2 + x3, data=mydata)  
summary(fit) # show results
```

```
# Other useful functions  
coefficients(fit) # model coefficients  
confint(fit, level=0.95) # CIs for model parameters  
fitted(fit) # predicted values  
residuals(fit) # residuals  
anova(fit) # anova table  
vcov(fit) # covariance matrix for model parameters  
influence(fit) # regression diagnostics
```

Diagnostické grafy

```
# diagnostic plots  
layout(matrix(c(1,2,3,4),2,2)) # optional 4 graphs/page  
plot(fit)
```



Porovnanie a Validácia modelov

```
# compare models  
fit1 <- lm(y ~ x1 + x2 + x3 + x4, data=mydata)  
fit2 <- lm(y ~ x1 + x2)  
anova(fit1, fit2)
```

```
# K-fold cross-validation  
library(DAAG)  
cv.lm(df=mydata, fit, m=3) # 3 fold cross-validation
```

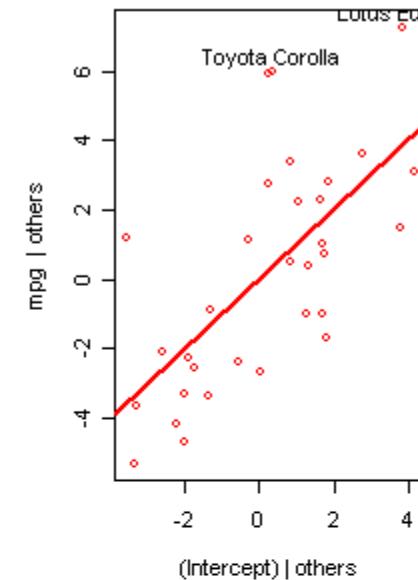
Regresná Diagnostika

```
# Assume that we are fitting a multiple linear regression  
# on the MTCARS data  
library(car)  
fit <- lm(mpg~disp+hp+wt+drat, data=mtcars)
```

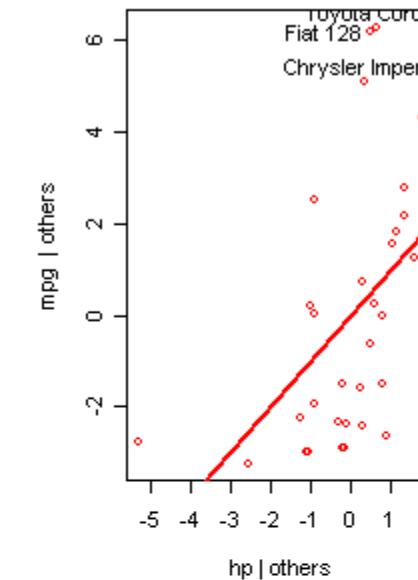
Regresná diagnostika

```
# Assessing Outliers  
outlierTest(fit) # Bonferonni p-value for most extreme obs  
qqPlot(fit, main="QQ Plot") #qq plot for studentized resid  
leveragePlots(fit) # leverage plots
```

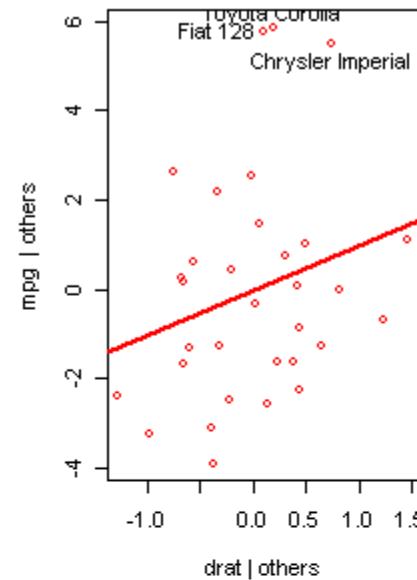
Leverage Plot



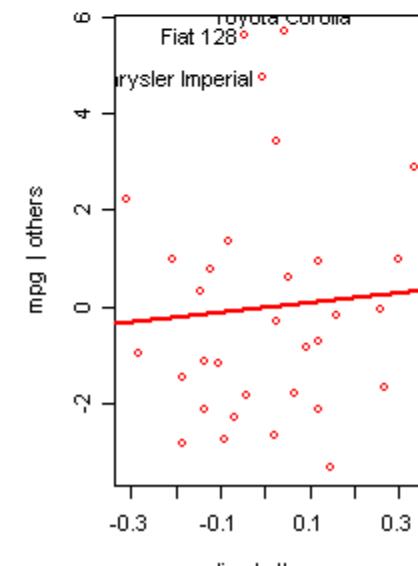
Leverage Plot



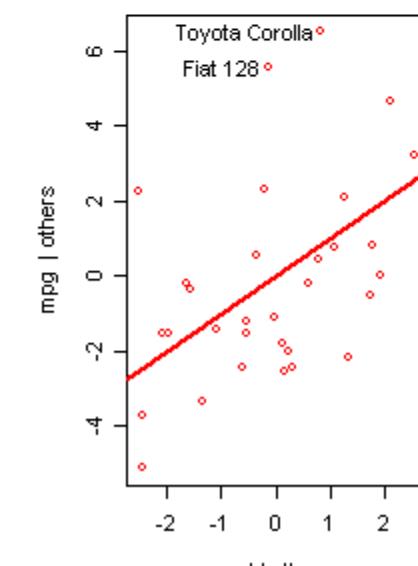
Leverage Plot



Leverage Plot

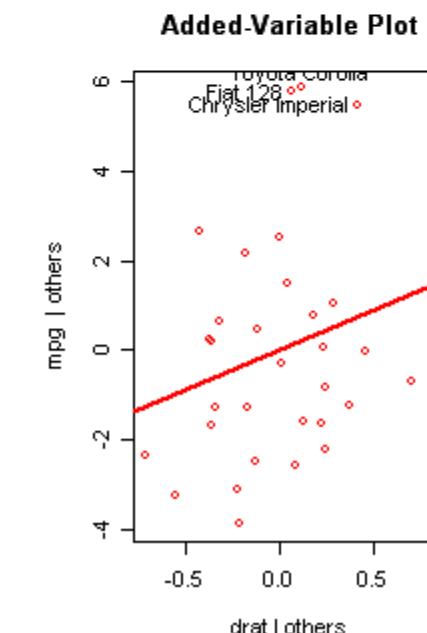
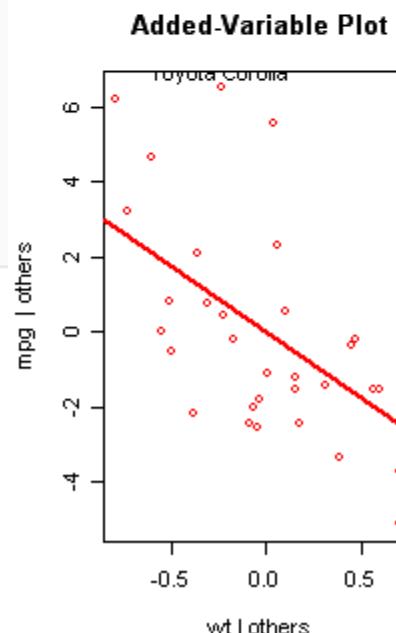
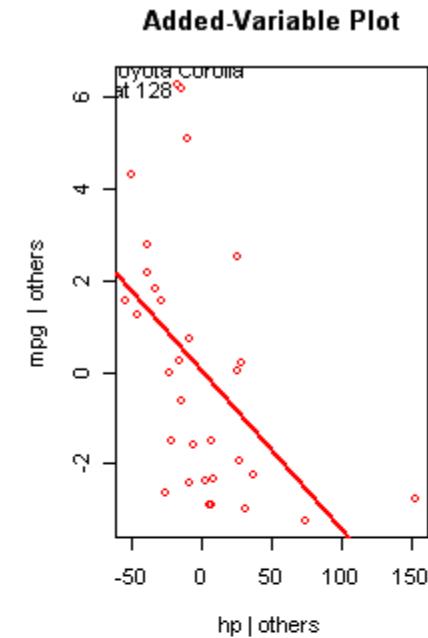
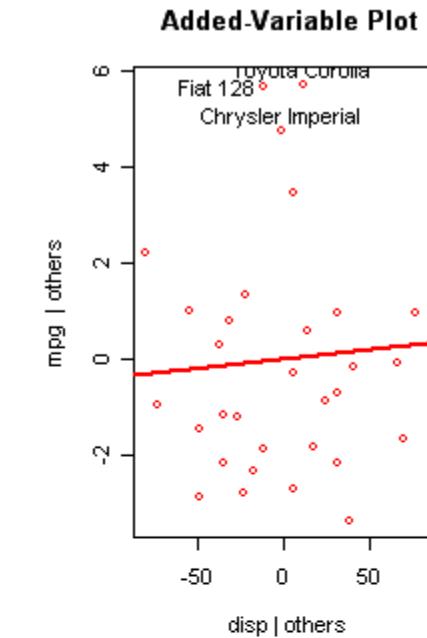
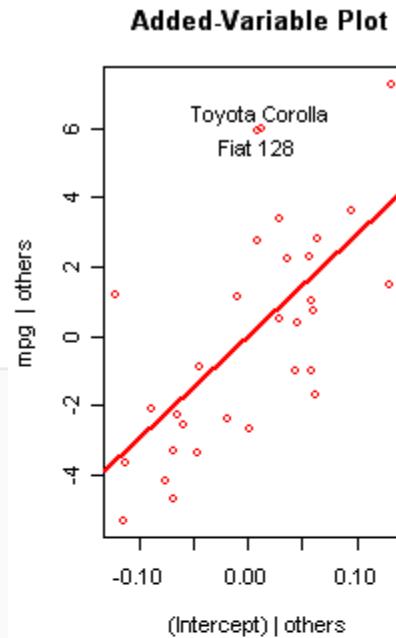


Leverage Plot

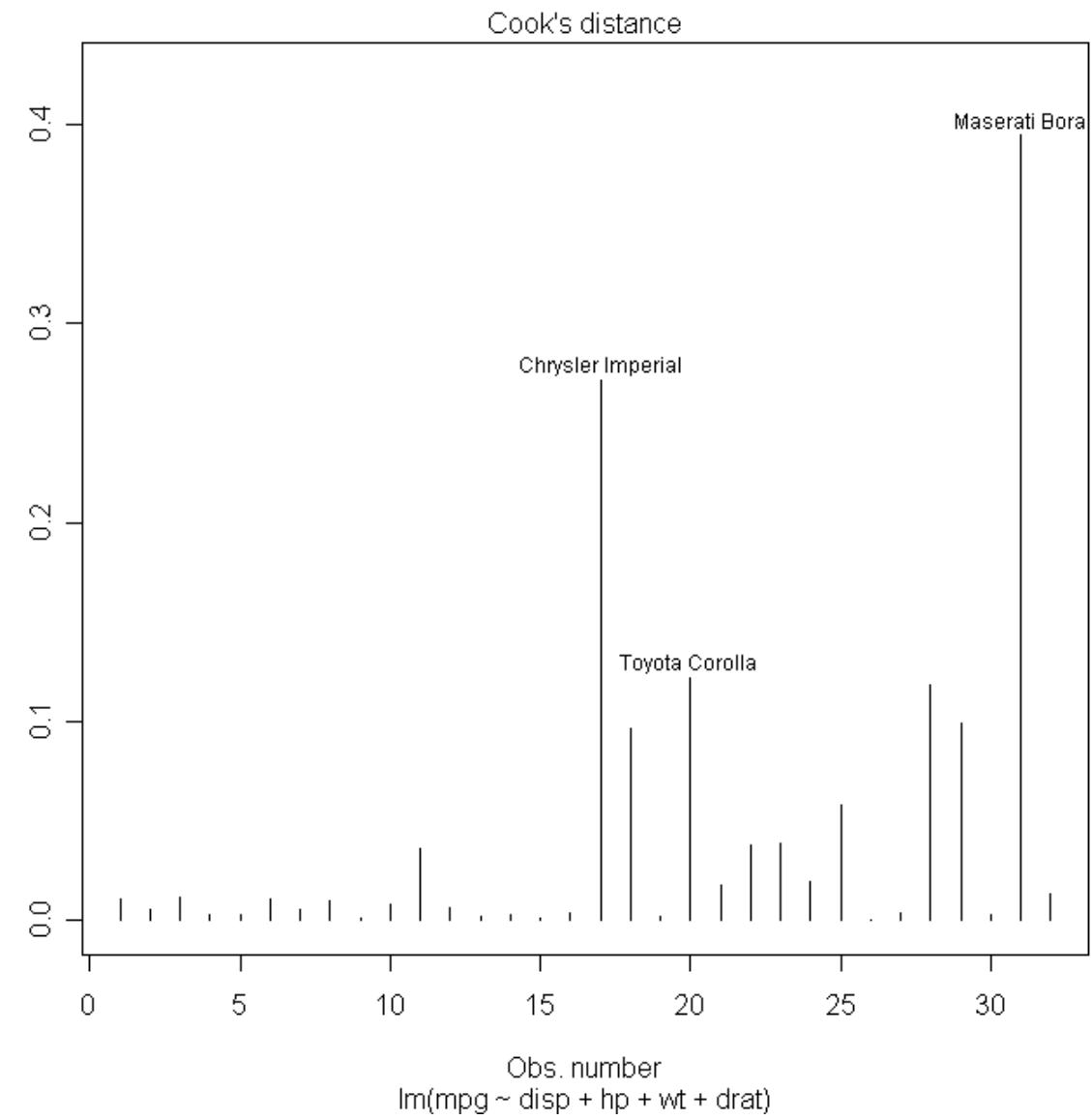
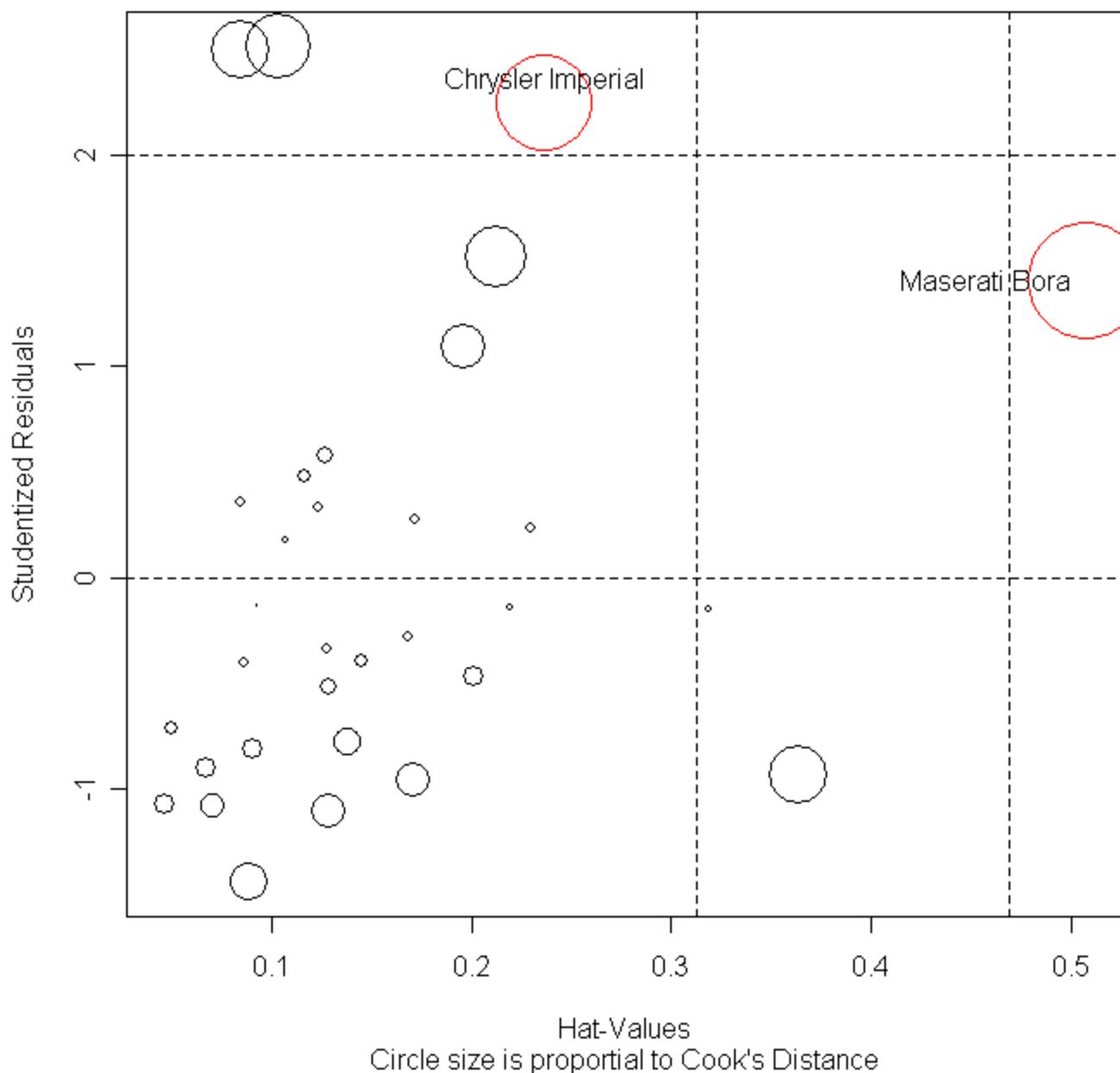


Vplyvné pozorovania

```
# Influential Observations  
# added variable plots  
av.Plots(fit)  
# Cook's D plot  
# identify D values > 4/(n-k-1)  
cutoff <- 4/((nrow(mtcars)-length(fit$coefficients)-2))  
plot(fit, which=4, cook.levels=cutoff)  
# Influence Plot  
influencePlot(fit, id.method="identify", main="Influence Plot",  
sub="Circle size is proportional to Cook's Distance" )
```



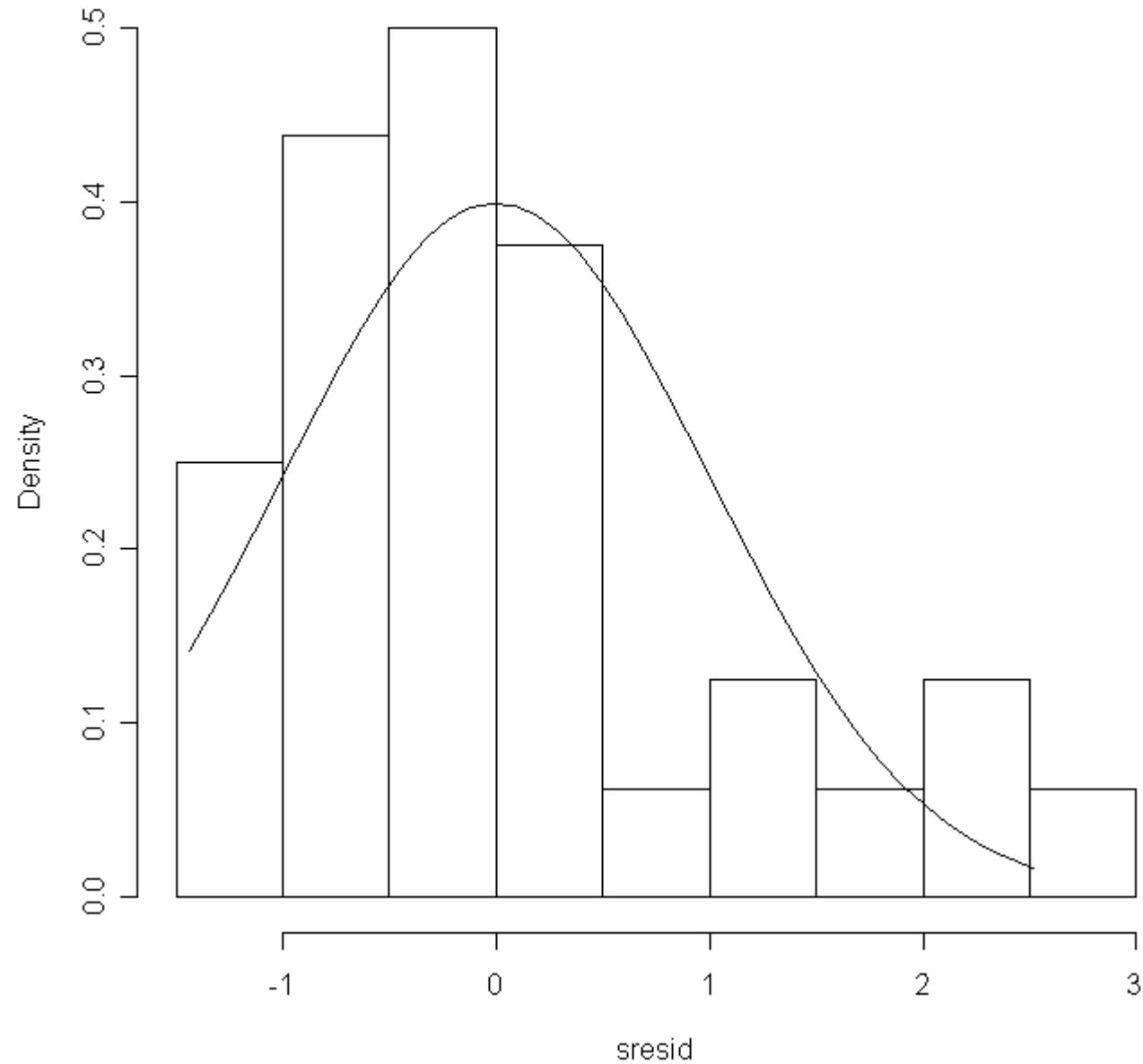
Influence Plot



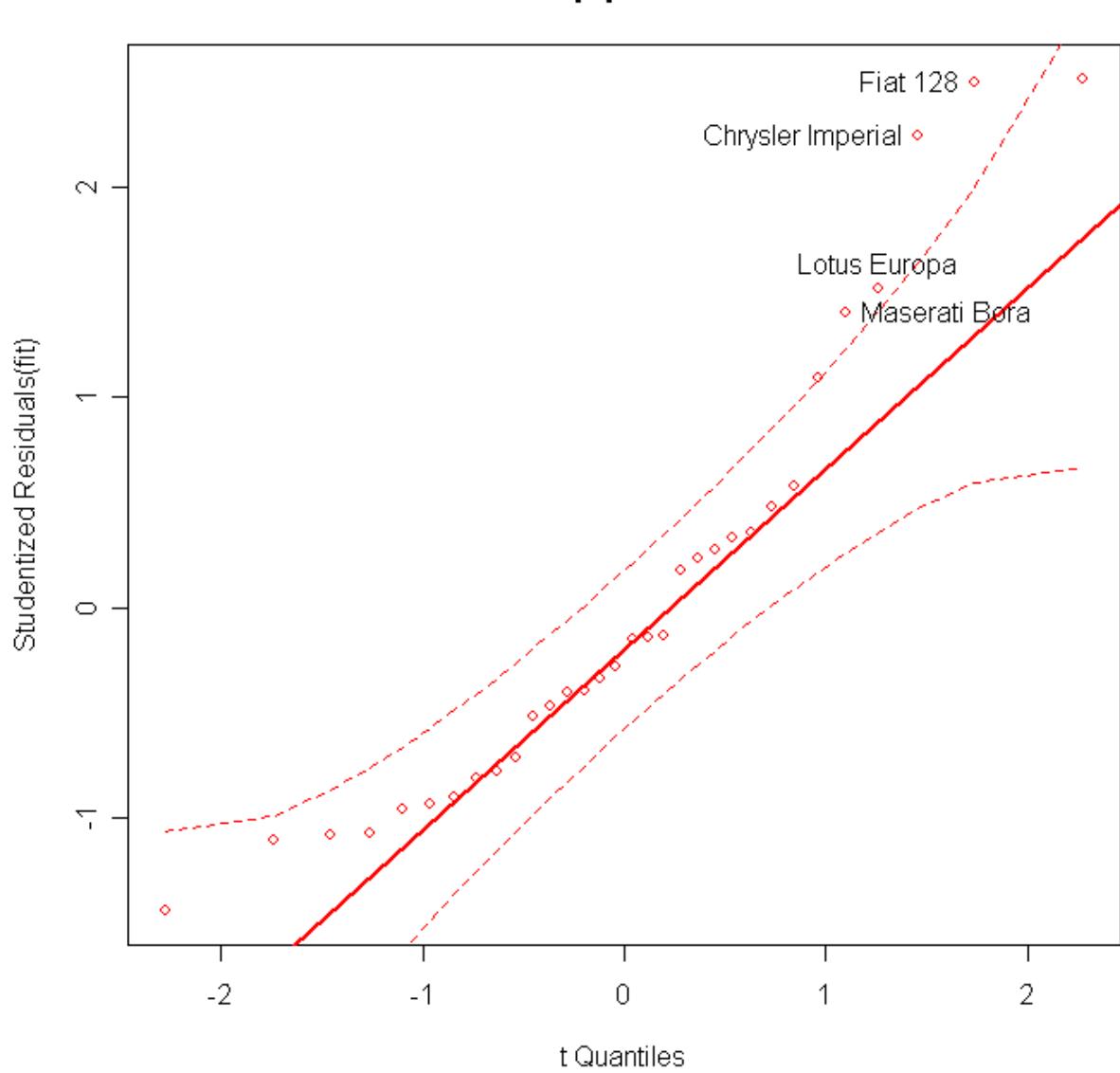
Non-normality

```
# Normality of Residuals  
# qq plot for studentized resid  
qqPlot(fit, main="QQ Plot")  
# distribution of studentized residuals  
library(MASS)  
sresid <- studres(fit)  
hist(sresid, freq=FALSE,  
     main="Distribution of Studentized Residuals")  
xfit<-seq(min(sresid),max(sresid),length=40)  
yfit<-dnorm(xfit)  
lines(xfit, yfit)
```

Distribution of Studentized Residuals



q-q Plot



Skúmanie závislostí medzi premennými (javmi, znakmi)

Symbolika:

- ✓ ZP – závislá premenná Y (resp. $Y=A$), resp. vysvetľovaná, modelovaná, cieľová, angl. dependent variable, response variable, target
- ✓ NZP – nezávislá premenná X, (resp. A, B,), resp. vysvetľujúca, angl. independent variable, input variable, exploration variable

1. Skúmanie závislostí medzi premennými (javmi, znakmi) - A a B

ZP	NZP	Názov závislosti	Štat. metóda/y	Jazyk R
A	B	asociácia	<ul style="list-style-type: none">Analýza kontingenčnej tabuľkyChi2 test závislosti (asociácie)	Describe – Table Analysis

2a. Skúmanie závislostí medzi premennými (javmi, znakmi) - Y a A

ZP	NZP	Názov záv.	Štat. metóda/y	Jazyk R
Y: Y_{A1} a Y_{A2}	A - faktorová premenná (Classification variable) len binárna: A_1 a A_2 $k = 2$		t-test (predtým ešte F-test): <ul style="list-style-type: none"> • pre nezávislé výbery • pre závislé výbery - párový 	ANOVA - t-test: <ul style="list-style-type: none"> • Two Sample (F-test je súčasť výstupu) • Paired

2b. Skúmanie závislostí medzi premennými (javmi, znakmi) – Y a A

ZP	NZP	Názov záv.	Štat. metóda/y	Jazyk R
Y: Y_{A_1} až Y_{A_k}	A – faktorová premenná (Classification variable) množná: A_1 až A_k $k > 2$		ANOVA: <ul style="list-style-type: none"> 1. a) jednofaktorová b) viacfaktorová 2. a) parametrická b) neparametrická 3. a) vyvážená b) nevyvážená 	ANOVA: <ul style="list-style-type: none"> 1. a) len 1-faktorová (One-Way ANOVA) <ul style="list-style-type: none"> b) 1 aj viacfaktorová (Linear Models) 2. a) parametrická (One-Way ANOVA, Linear Models) <ul style="list-style-type: none"> b) neparametrická (Nonparametric One-Way ANOVA) 3. a) vyvážená (One-Way ANOVA) <ul style="list-style-type: none"> b) nevyvážená (Linear Models)

3a. Skúmanie závislostí medzi premennými (javmi, znakmi) - Y a X

ZP	NZP	Názov záv.	Štat. metóda/y	Jazyk R
Y	X	Korelácia jednoduchá	<p>1. výpočet mier závislosti (korelácie):</p> <ul style="list-style-type: none"> • $\text{cov}(x,y) \in (-\infty, \infty)$ • $r(x,y) \in (-1, 1)$ <p>2. odhad regresného modelu jednoduchého:</p> <ul style="list-style-type: none"> • lineárneho alebo • nelineárneho 	<p>1. Multivariate – Correlations</p> <p>2. Regression:</p> <ul style="list-style-type: none"> • Linear Regression • Nonlinear Regression

3b. Skúmanie závislostí medzi premennými (javmi, znakmi) – Y a X

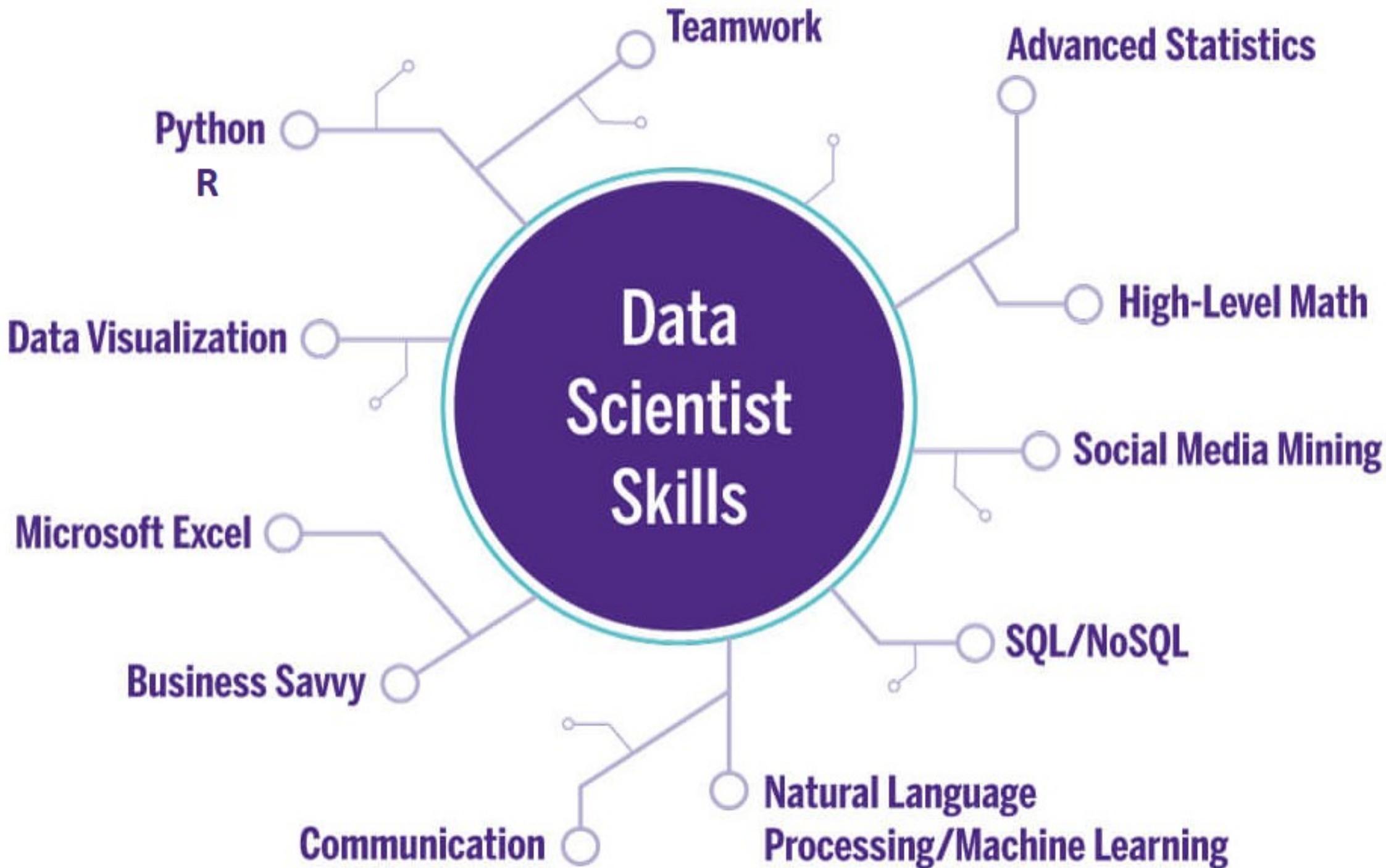
ZP	NZP	Názov záv.	Štat. metóda/y	Jazyk R
Y	$X_1, X_2, X_3, \dots, X_k$	Korelácia viacnásobná	<p>1. výpočet mier závislosti (korelácie):</p> <ul style="list-style-type: none"> kovariančná matica (S) korelačná matica (R) <p>2. odhad regresného modelu viacnásobného:</p> <ul style="list-style-type: none"> lineárneho alebo nelineárneho 	<p>1. Multivariate – Correlations</p> <p>2. Regression:</p> <ul style="list-style-type: none"> Linear Regression Nonlinear Regression

4. Skúmanie závislostí medzi premennými (javmi, znakmi) – Y a t (čas)

ZP	NZP	Názov záv.	Štat. metóda/y	Jazyk R
Y	t (čas)		Analýza časových radov	Time Series Analysis

5. Skúmanie závislostí medzi premennými (javmi, znakmi) – $Y=A$ a X_i, B_j

ZP	NZP	Názov záv.	Štat. metóda/y	Jazyk R
$Y=A$ $P(Y=A_1/X_i, B_j)$	$X_1, X_2, X_3, \dots, X_k$ $B_1, B_2, B_3, \dots, B_q$		Logistická regresia	Regression – Logistic Regression



Zabudované Datasety

The screenshot shows the RStudio interface. The top panel is the 'R data sets' browser, displaying a list of datasets from the 'datasets' package. The bottom panel is the 'Console' tab, showing the command `> data()` entered by the user.

Data set	Description
AirPassengers	Monthly Airline Passenger Numbers 1949-1960
BJSales	Sales Data with Leading Indicator
BJSales.lead (BJSales)	Sales Data with Leading Indicator
BOD	Biochemical Oxygen Demand
CO2	Carbon Dioxide Uptake in Grass Plants
ChickWeight	Weight versus age of chicks on different diets
DNase	Elisa assay of DNase
EuStockMarkets	Daily Closing Prices of Major European Stock Indices, 1991-1998
Formaldehyde	Determination of Formaldehyde
HairEyeColor	Hair and Eye Color of Statistics Students
Harman23.cor	Harman Example 2.3
Harman74.cor	Harman Example 7.4
Indometh	Pharmacokinetics of Indomethacin
InsectSprays	Effectiveness of Insect Sprays
JohnsonJohnson	Quarterly Earnings per Johnson & Johnson Share
LakeHuron	Level of Lake Huron 1875-1972
LifeCycleSavings	Intercountry Life-Cycle Savings Data
Loblolly	Growth of Loblolly pine trees
Nile	Flow of the River Nile

```
> data()
```

Data()

Apropos a Demo

- Príkaz **apropos ("nazov")** zobrazíme príbuzné príkazy k príkazu názov, pre príklady k nejakej téme použijeme príkaz example (tema)
- K dispozícii sú tiež **PDF manuály**, ktoré obsahujú základné manuály a referencie o funkciách
- Príkaz **demo()** vypíše zoznam všetkých dostupných demonštrácií.

The screenshot shows the RGui (64-bit) interface. In the R Console window, the following R code is displayed:

```
> apropos("lq")
[1] "evalq"      "evalqOnLoad"
> apropos("ls")
[1] ".colSums"    ".getKlevels"   ".signalSimpleWarning"
[4] "ar.ols"       "co.intervals" "colSums"
[7] "de.ncols"    "densCols"     "droplevels"
[10] "droplevels.data.frame" "droplevels.factor" "evalSource"
[13] "formals"      "formals<-"    "getAllSuperClasses"
[16] "getLoadedDLLs" "ifelse"       "labels"
[19] "labels.default" "levels"       "levels.default"
[22] "levels<-"     "levels<-factor" "limitedLabels"
[25] "ls"           "ls.diag"      "ls.print"
[28] "ls.str"       "lsf.str"      "lsfit"
[31] "nlevels"      "nls"          "nls.control"
[34] "NLSstAsymptotic" "NLSstClosestX" "NLSstLfAsymptote"
[37] "NLSstRtAsymptote" "occupationalStatus" "OlsonNames"
[40] "residuals"    "residuals.glm" "residuals.lm"
[43] "restartFormals" "symbols"      "sys.calls"
[46] "weighted.residuals"
> demo()
> demo(colors)
```

Below the console, the command `demo(colors)` is entered. To the right of the console, a color palette titled "Click or hit ENTER for next page" is shown with three colored squares: deepskyblue, turquoise1, and cyan.

Zlúčenie/Merging Údajov

```
# merge two data frames by ID  
total <- merge(data frameA,data frameB,by="ID")
```

```
# merge two data frames by ID and Country  
total <- merge(data frameA,data frameB,by=c("ID", "Country"))
```

```
total <- rbind(data frameA, data frameB)
```

Čo sa Oplatí Prečítať?

Slovensko a česko

- Albatrosmedia
- Kopp
- Grada
- Wolters Kluwer
- BEN
- Veda

Zahraničie

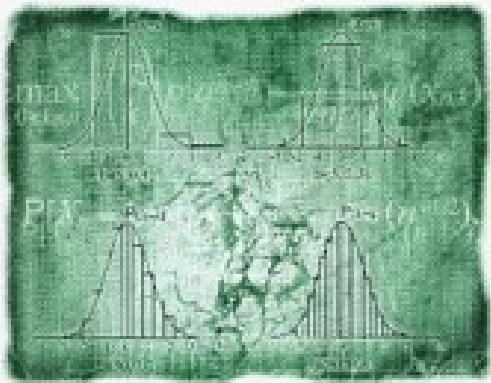
- O'Reilly
- Manning
- Packt
- Apress
- Wiley
- No Starch Press

YouTube Tutoriály

- [IT Academy](#)

KAREL JUHA a JOSEF ŠTĚPÁN

Pravděpodobnost
a matematická
statistika



matfyzpress
www.matfyz.cz
e-mail: info@matfyz.cz
ISBN 978-80-7372-200-7

Průvodce základními statistickými metodami

Marek Budík
Marek Křížek
Bohumil Macek



- Ačkoli se vyučují v dívčích místnostech pouze v matematice a fyzice
- Ačkoli je statistika matematickou vědou
- Pouze významnou výzvou je vyučování
- Matematika je vyučována pomocí expozic
- Ačkoli je často vyučována na základě výuky fyziky
- Matematika je vyučována pomocí expozic

GRADA

Stanislav Šípek

Statistiká

bez předchozích
znalostí

Průvodce
pro samouky

Vymělné terminologie
a myšlenky v praxi
Nezávadného hledání
bez složitých teorií
Kombinace statistiky
a výpočetní techniky
Propojení i využití vzdálené
do téma statistiky

OPRESS
Osborne

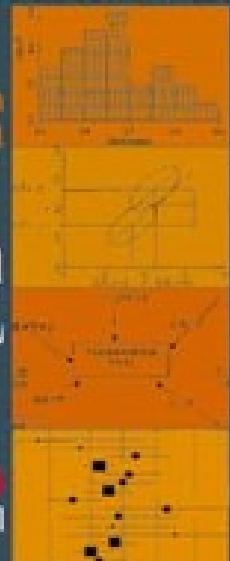


Přehled
STATISTICKÝCH METOD
zpracování dat

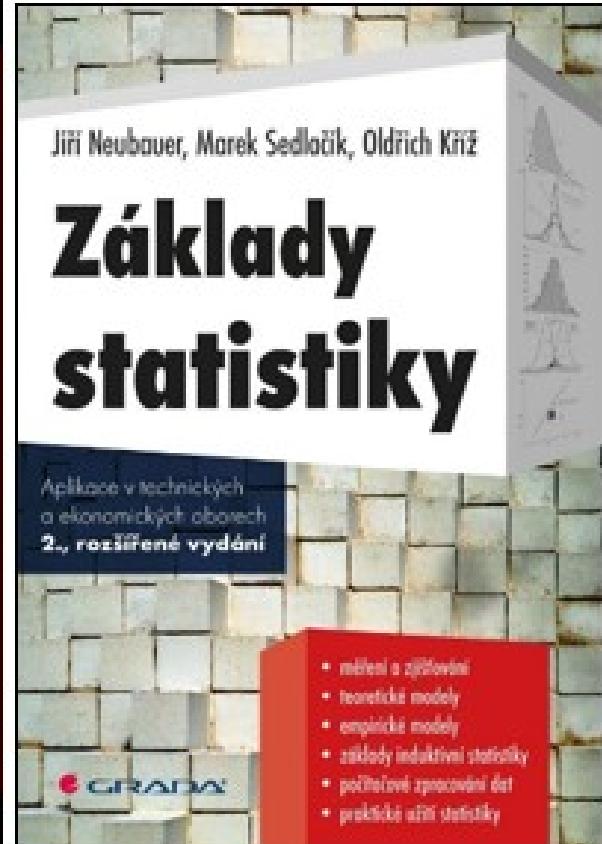
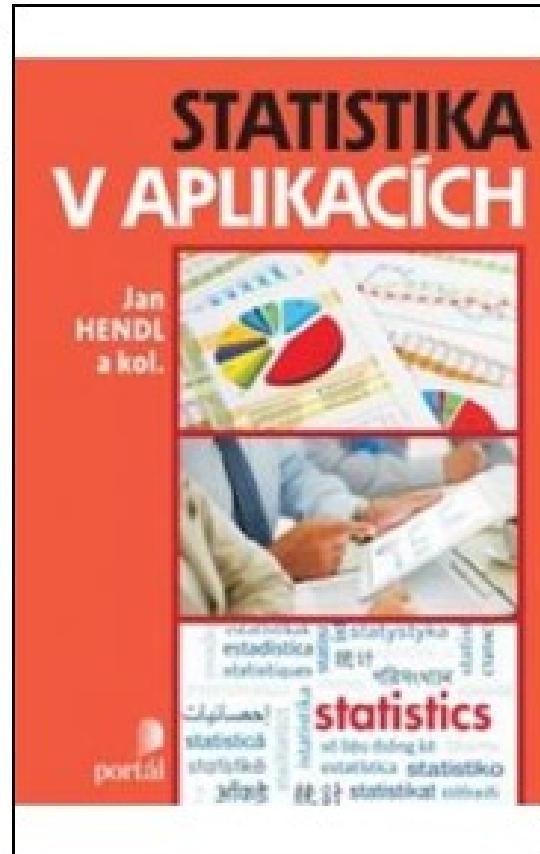
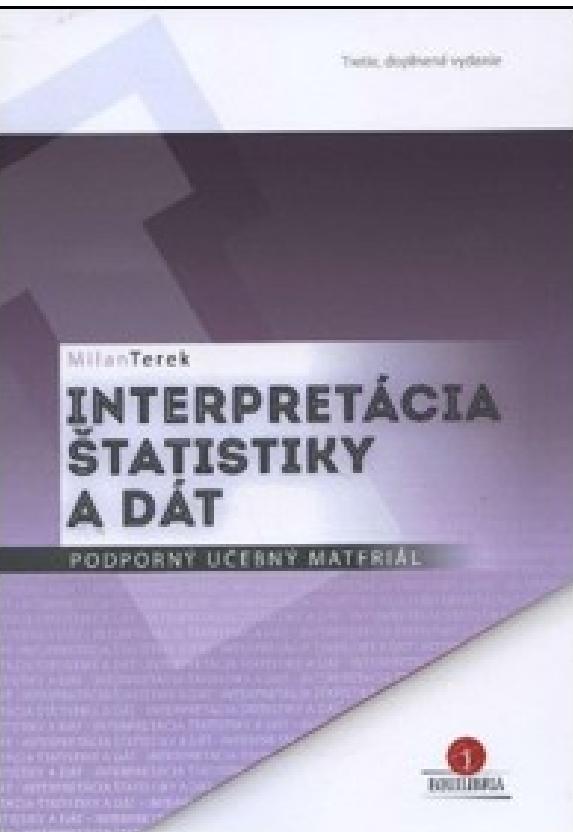
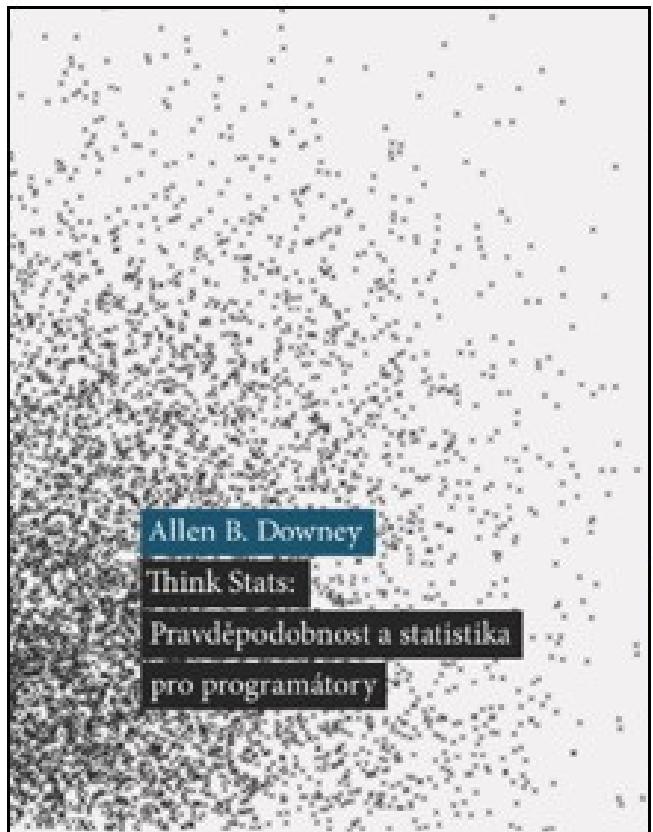
Analýza
a metranoanalýza dat

Jan
HENDL

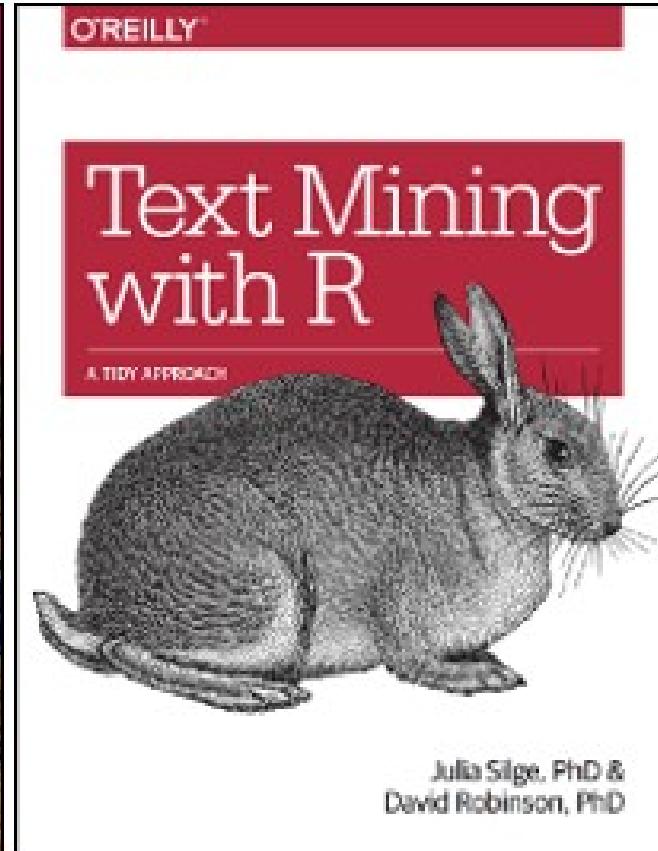
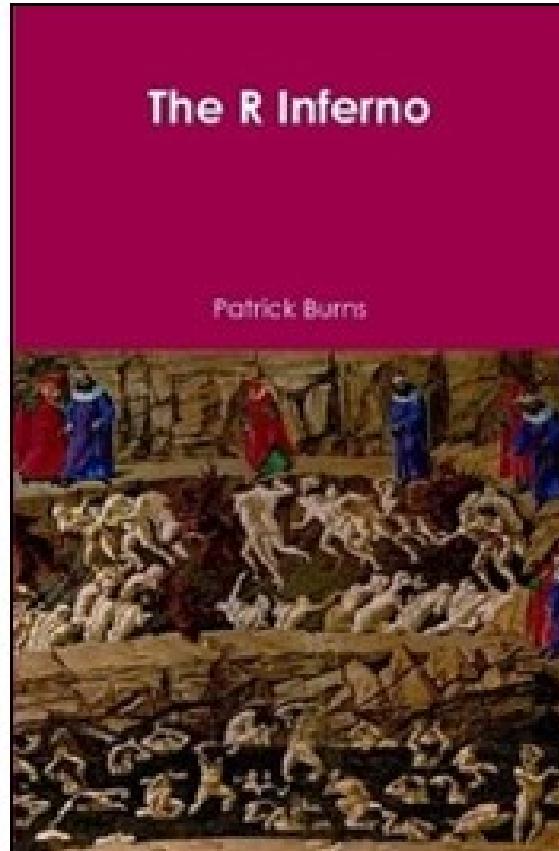
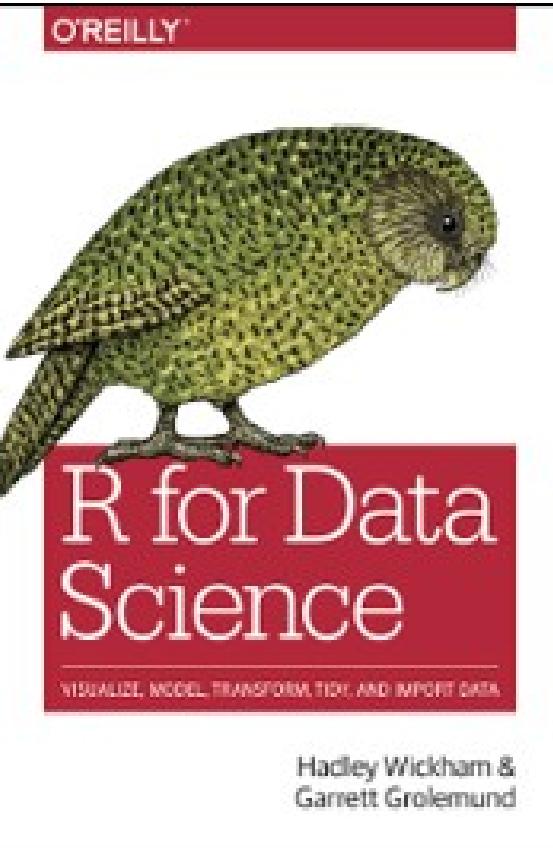
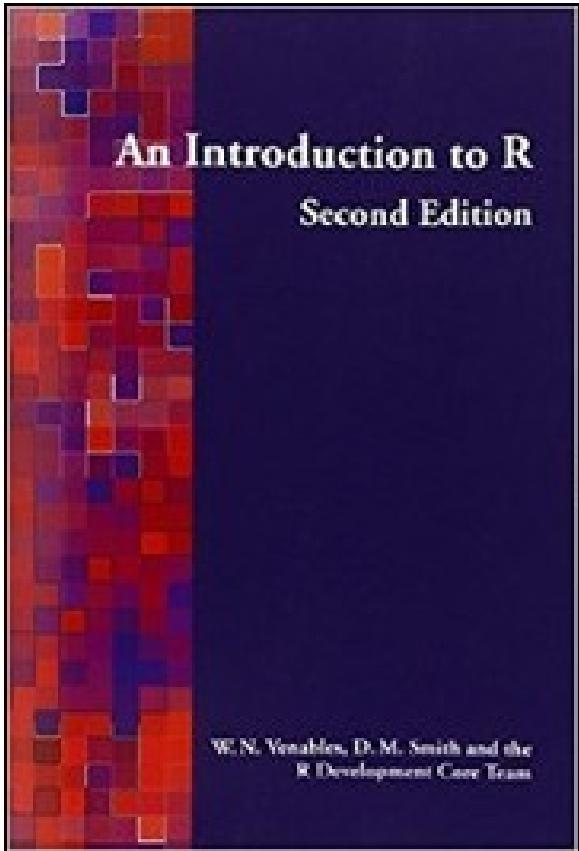
portál



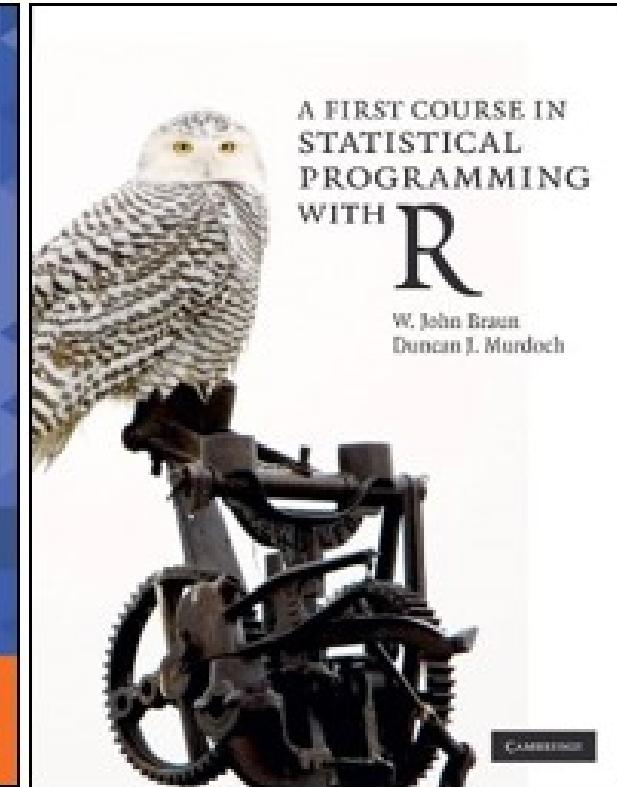
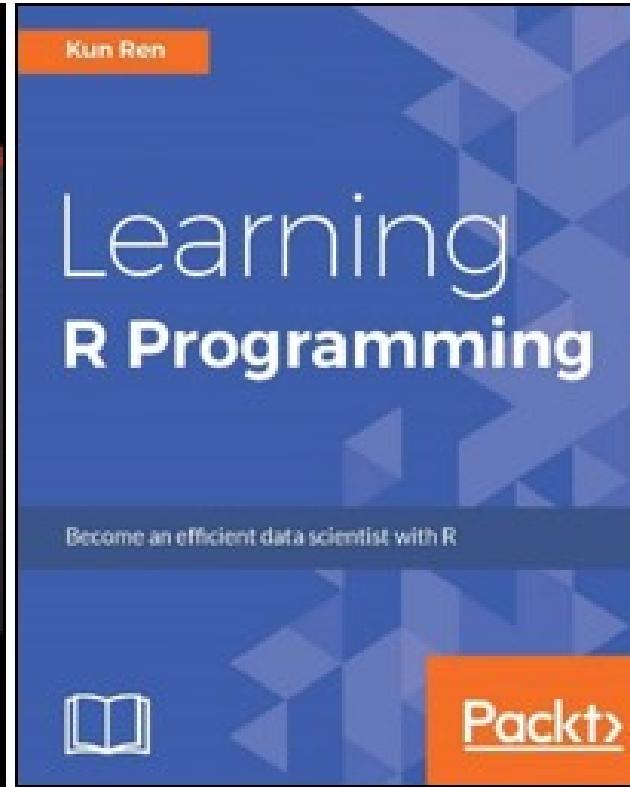
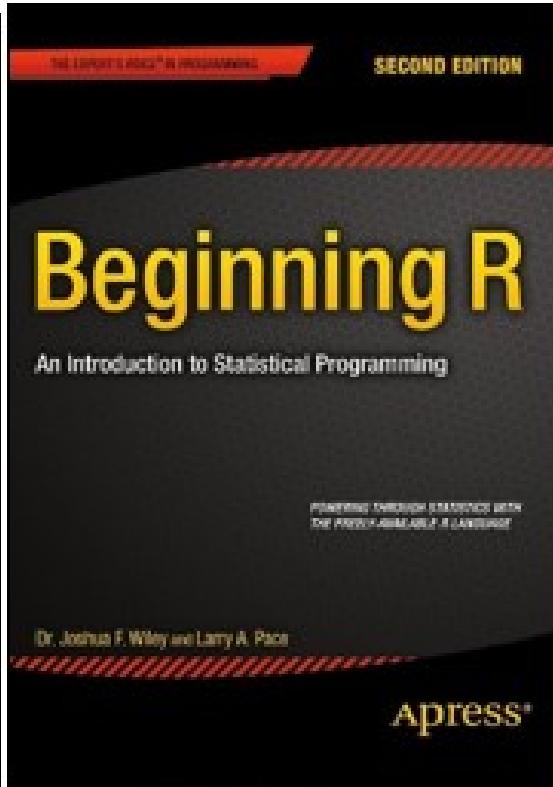
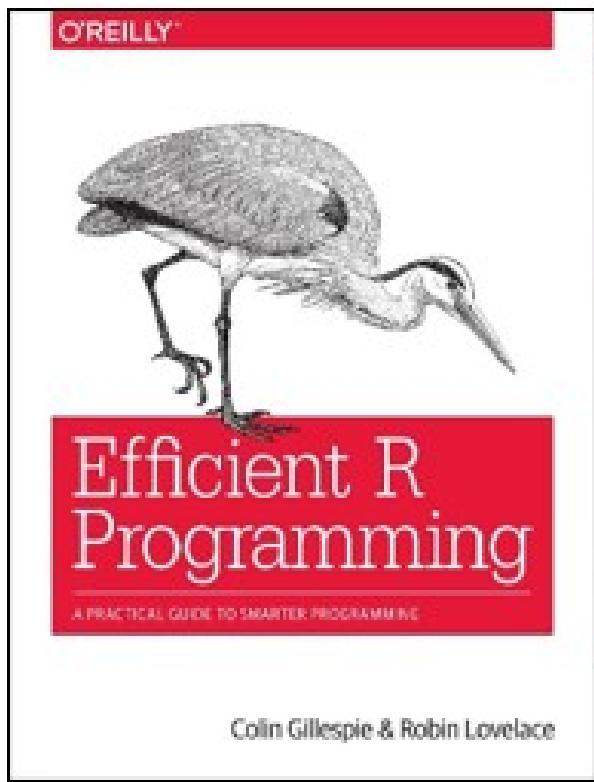
Čo sa oplatí/neoplatí prečítať SK/CZ?



Čo sa oplatí/neoplatí prečítať SK/CZ?



Čo sa oplatí/neoplatí prečítať EN



Čo sa oplatí/neoplatí prečítať EN

Home

PUBLIC

 Questions**Tags**

Users

Companies

COLLECTIVES

 Explore Collectives

TEAMS

 Create free Team

Looking for your Teams?

Tags

A tag is a keyword or label that categorizes your question with other, similar questions. Using the right tags makes it easier for others to find and answer your question.

[Show all tag synonyms](#)

r

Popular

Name

New



r

489303 questions 113 asked today, 800 this week

 javascript

2493694 questions 401 asked today, 2384 this week

 android

1403892 questions 159 asked today, 856 this week

 jquery

1033870 questions 46 asked today, 229 this week

 reactjs

454618 questions 219 asked today, 1319 this week

 arrays

412869 questions 34 asked today, 227 this week

 ruby-on-rails

336768 questions 16 asked today, 102 this week

 sql-server

329798 questions 45 asked today, 255 this week

 angular

295326 questions 72 asked today, 444 this week

 angularjs

262780 questions 7 asked today, 27 this week

 regex

257806 questions 25 asked today, 162 this week

 ruby

227846 questions 11 asked today, 77 this week

Dáta, analýzy a štatistiky

Verejná skupina · 1,5 tis. členov

+ Pozvati

Informácie

Diskusia

Členovia

Podujatia

Médiá

Súbory



Create a public post...



Fotka/video



Označiť ľudí



Pocit/aktivita

Nová aktivity ▾



Pavol Skapik zdieľa odkaz.

Administrátor · 29. januára o 19:18 ·

...

Marek Krajčí vyhlásil, že po dopočítaní oneskorených údajov sa podiel pozitívnych testov v Bratislave dostał až na 1,3 percenta

BSK: Výsledky testovania v mestských častiach Bratislavы hovoria o podiele pozitívnych 0,81%, resp. všetko čo je k dispozícii je 0,88.

Čím si prosím vysvetľujete rozdiel? Ďakujem... Zobrazí viac

Výsledky plošného skríningu



Informácie

Skupina je len súkromným projektom, v žiadnom prípade zverejnené názory a príspevky nie sú oficiálnymi výstupmi.

Verejná

Členov skupiny a ich príspevky bude vidieť ktokoľvek.

Viditeľná

Túto skupinu nájde ktokoľvek

Bratislava, Slovakia

Skupina o Všeobecné

Populárne témy v príspevkoch

census (1)

demograf...





Hľadať na Facebooku



Miroslav



Data Analysts, Data Engineers & Data Scientists - Czech&Slovak Group

Verejná skupina · 2,3 tis. členov



+ Pozvat

Informácie

Diskusia

Oznámenia

Miestnosti

Témy

Členovia

Podujatia

Médiá



Create a public post...



Fotka/video



Označiť ľudí



Pocit/aktivita

Oznámenie · 1

Zobrazit všetko



Vojta Roček zdieľa odkaz.

★ Administrátor · 21. apríla 2017 · Prague, Praha, Česká republika ·

...

Zajímavé veci, na ktere jste prisli a chcete je posdílet. Do popisu klidne piste, proc vam to pripada zajímavé. Otazky, odpovedi, rady. Přidejte další datalidi, a ptejte se na co chcete! (Ale až po prostudování <https://www.hash.cz/inferno/otazky.html>).

Nabídky a poptávky po práci jsou v sesterské skupině -
<https://www.facebook.com/groups/1788236724824404/>

Pokud se chystáte místních 1600 lidí oslovit zdarma s nabídkou na svou komerční akci (spam), tak to udelejte to dobre (...). [Zobrazit viac](#)

[Zobrazit preklad](#)

Informácie

Dataflow.cz Dataflow

Twitter: <https://twitter.com/dataflowcz>

Web: <http://dataflow.cz...> [Zobrazit viac](#)

● Verejná

Členovia skupiny a ich príspevky bude vidieť ktokoľvek.

● Viditeľná

Túto skupinu nájde ktokoľvek

● Prague, Czech Republic

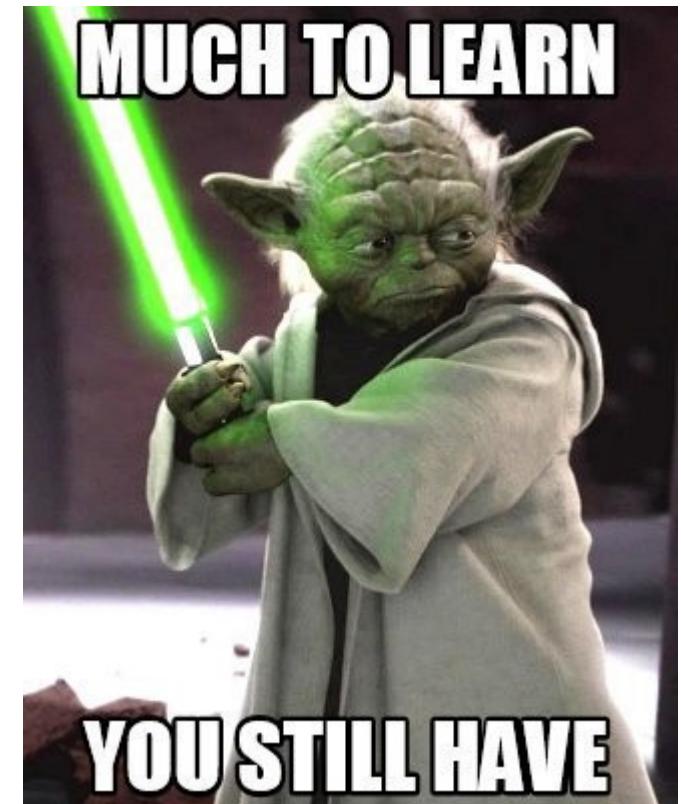
● Skupina o Všeobecné

Nedáve médiá



Čo ti odporúčam si pozrieť?

1. <https://google.github.io/styleguide/Rguide.xml>
2. <https://www.statmethods.net/input/exportingdata.html>
3. http://www.sr.bham.ac.uk/~ajrs/R/r-function_list.html
4. <https://www.statmethods.net/management/index.html>
5. <https://cran.r-project.org/web/packages/stringr/vignettes/regular-expressions.html>



Šup do záložiek

Najdôležitejšie Klávesové Skratky

Práca s IDE

- Ctrl + D Delete zmaž riadok
- **Ctrl + Space** Asistent kódu
- **Ctrl + /** Komentáre
- Ctrl + A Označ všetko
- **Alt + /** Dokonči slovo
- Ctrl + F Hľadanie a náhrady
- Ctrl + Shift + F Kompakt režim
- Ctrl + Shift + S Ulož všetko

Práca s Browserom

- Ctrl + T Vytvor nový tab
- Ctrl + W Zatvor aktuálny tab
- Ctrl + Shift + W Zatvor všetky taby
- **Ctrl + Shift + T** Otvor posledný tab
- Ctrl + Shift + J/F12 Web console
- **F11** Fullscreen

F5 nie je Spustenie, ale Refresh

Efektívne Používanie Klávesnice



Špeciálne znaky, kde ich nájst' na klávesnici



Operátory

+ Sčítavanie, Spájanie
* Násobenie, Opakovanie
- Odčítanie
/ Delenie
% Zvyšok, modulo
@ Dekorátor, Nás. matic

Porovnávanie

< > Väčšie, Menšie
= Rovnosť, Priradenie
! Nerovnosť
" Úvodzovky
\ Špeciálne znaky

Oddelovače

, Prvkov
. Atribútov
; Blokov, Klúčov
: Príkazov
Poznámky
Komentár

Bitové operácie

& Prienik, A, AND
| Zjednotenie, OR
^ XOR
~ Inverzie

Zátvorky

() Zátvorky, Volanie
{ } Slovníky, Formát
[] Zoznamy, Indexy
Ostatné
Súčasť mena
\$ Nevyužité



Ako sa s nami spojiť?



ADRESA: IT Academy, s. r. o.

Budova KOLOSEO prízemie
Tomášikova 50/A
831 04 Bratislava



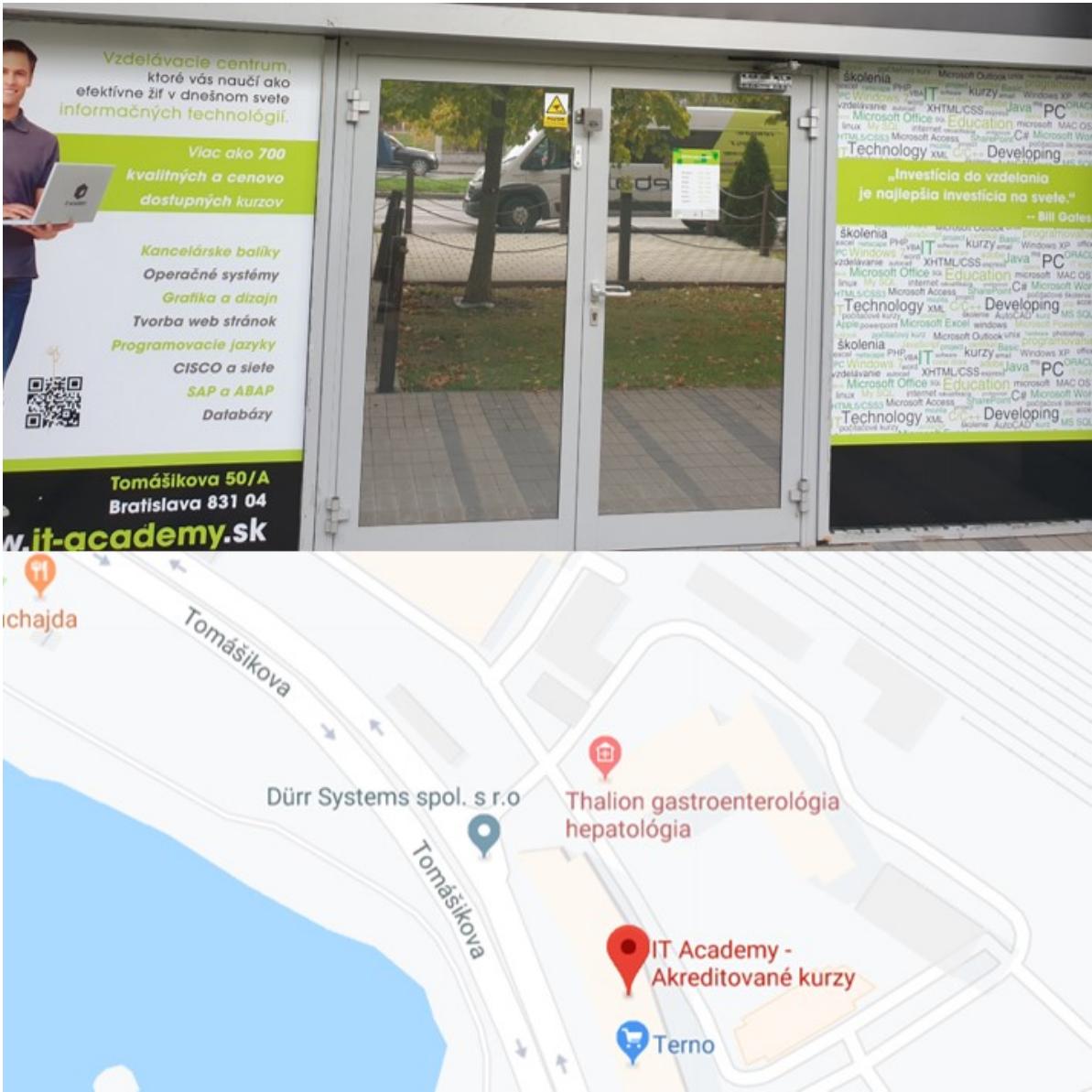
WEB: www.it-academy.sk



E-MAIL: info@it-academy.sk



TELEFÓN: +421 917 095 406



Ako vieme pomôc?

#Copywriting

#Školenia

#Zamestnanci

#Pomáhame

#Rast

#Projekty

#Certifikácie

#Kurzy

#Tréningy

#Vzdelávanie

#PPC Kampane

#Elearning

#Mentoring

#Konzultácie

#Online

#Programovanie

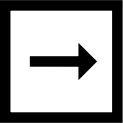
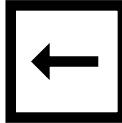
#Vývoj

#Marketing

#Reklama

#Prenájom Techniky

Mrkni na náš YouTube kanál a daj odber

 [WWW.YOUTUBE.COM/C/IT-ACADEMYSK](https://www.youtube.com/c/IT-ACADEMYSK) 

Vyber si online kurz

Nauč sa programovať, tvoriť webstránky a grafiku, manažovať alebo sa zameraj na osobný rozvoj. Všetko jednoducho vďaka našim online kurzom z pohodlia tvojho domova.

Ročné predplatné na
všetky online kurzy

~~2299.99€~~

299.99€

Prístup pre Teba do všetkých
aktuálnych aj pripravovaných
online kurzov

12 mesačná platnosť

Kúpiť teraz

Zadarmo

- A. Ďalšie kurzy NSPC
- B. Grow with Google

1. NESTRAŤ PRÁCU (Zamestnaní, SZČO,
na materskej)
2. KOMPAS/REPAS (Nezamestnaní)
3. YouTube Kanál (IT Academy)

Platené

1. Moje kurzy na www.vita.sk