

## R pro analýzu vědeckých dat

Mgr. Lucie Hošková







TCS





#### Obsah

- Instalace R a Rstudio
- R vs. Rstudio: jaký je rozdíl?
- Stáhnutí dat
- Nastavení vizuálu
- Vytvoření projektu
- Popis oken
- Klávesové zkratky





#### Instalace R a RStudio

R

### Rstudio

https://mirrors.nic.cz/R/

Desktop verze: odkaz

dostupné pro Windows, macOS i Linux dostupné pro Windows, macOS i Linux



#### R vs. Rstudio: jaký je rozdíl?

#### R

• open source jazyk pro zpracovávání dat

#### RStudio

- open source vývojové prostředí pro práci s jazykem R
- zobrazuje kromě zdrojového řádku navíc oddělenou konzoli či globální prostředí, celkovo uživatelsky přívětivější

#### Budeme pracovat pouze přímo uvnitř prostředí programu RStudio!

National Czech Programme

#### Stáhnutí dat

 Během školení se bude pracovat na předem připravených datech. Stáhněte si proto prosím následující soubor "newborns.txt", a uložte si jej do nové složky, která bude sloužit pouze pro soubory k tomuto školení. Neukládejte do ní prosím v tuto chvíli žádné další soubory a uložte si ji na svém počítači na takovém místě, kde ji budete schopni opět snadno najít.





#### Nastavení vizuálu RStudia

- Předtím než začnete v Rstudiu pracovat, je dobré nastavit si vizuál rozhraní.
- Můžete samozřejmě pracovat i se základním nastavením, ale doporučuji si i tak ale nastavení alespoň prohlédnout.
- Na následujících snímcích máte návod krok za krokem, jak tuto změnu vizuálu provést.

<pre>O I I Constant I I I I I I I I I I I I I I I I I I I</pre>	File Edit Code View Plots Session Build Debug Profile Tools	
Console Terminal & Badground Jobs « R R432:-/** R version 4.3.2 (2023-10-31 ucrt) "Eye Holes" Copyright (c) 2023 The R Foundation for Statistical Computing Platform: x86_64-w64-mingw32/x64 (64-bit) R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details. R is a collaborative project with many contributors. Type 'contributors()' for more information and 'citation()' on how to cite R or R packages in publications. Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'q()' to quit R. >	• ♥ ♥	8
R version 4.3.2 (2023-10-31 ucrt) "Eye Holes" Copyright (C) 2023 The R Foundation for Statistical Computing Platform: x86_64-w64-mingw32/x64 (64-bit) R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details. R is a collaborative project with many contributors. Type 'contributors()' for more information and 'citation()' on how to cite R or R packages in publications. Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'q()' to quit R. >	Console     Terminal ×     Background Jobs ×       R     R4.3.2 · ~/	Environment History Connections Tutorial
Type 'contributors()' for more information and         'citation()' on how to cite R or R packages in publications.         Type 'demo()' for some demos, 'help()' for on-line help, or         'help.start()' for an HTML browser interface to help.         Type 'q()' to quit R.         >	R version 4.3.2 (2023-10-31 ucrt) "Eye Holes" Copyright (C) 2023 The R Foundation for Statistical Computing Platform: x86_64-w64-mingw32/x64 (64-bit) R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details.	Environment is empty
🗌 🧮 Zoom	<pre>R is a collaborative project with many contributors. Type 'contributors()' for more information and 'citation()' on how to cite R or R packages in publications. Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'q()' to quit R. &gt;</pre>	Files       Plots       Packages       Help       Viewer       Presentation         Image: Second state of the
		Zoom



#### National Czech Programme

RStudio

#### meosc

– 0 X



meosc



- Všechny změny které provádíte můžete vidět na příkladu kódu v pravé části okna.
- Editor font size: změní velikost textu kódu
- Help panel size: změní velikost textu nápovědy
- Editor theme: určuje styl a jeho barevné rozložení kódu
  - Doporučuji vybrat styl, který jasně barevně odlišuje čísla, poznámky, funkce a text (např. Cobalt).



### Vytvoření projektu

- Jako první krok pro přípravu prostředí si vytvoříme Projekt.
- Projekt je "složka", do které se nám budou ukládat veškeré naše skripty, objekty, načtená data a historie naší práce.
- RStudio můžete samozřejmě používat i bez projektů, zvlášť když si jen rychle potřebujete ve skriptu projet nějaký výpočet. Pro zpracovávání výzkumných dat v jejich celku ovšem doporučuji projekty používat.
- Vytvořením projektu vás provedou následujících několik snímků.

ramme Czech				
RStudio				- 0 ×
File Gode View Blots Generate Build Debug Profile Tools Help				R Project: (None)
Console Terminal × Background Jobs ×	-	Environment History Connections Tut	orial	
R 4.3.2 · ~/ ≈		🕣 🕞 📅 Import Dataset 🔹 🔮 86 MiB 🗸	1	≣ List •   @ •
		R 👻 🛑 Global Environment 👻		Q
R version 4.3.2 (2023-10-31 ucrt) "Eye Holes" Copyright (C) 2023 The R Foundation for Statistical Computing Platform: x86_64-w64-mingw32/x64 (64-bit) R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details.	Environment is empty			
R is a collaborative project with many contributors. Type 'contributors()' for more information and		Files Plots Packages Help Viewer	Presentation	_
'citation()' on how to cite R or R packages in publications.		Sew Folder Sew Blank File - Sev De	elete 📑 Rename	🍓 More 🔹 🕓
Type 'dema()' for some demos 'help()' for $e_{1}$ ine help or		🗌 🏠 Home		
help.start()' for an HTML browser interface to help.		A Name	Size	Modified
Type $'a()'$ to quit R.		. Rhistory	146 B	Apr 2, 2025, 11:45 AM
		SRECYCLE.BIN	400 B	F-F 4 2024 2:40 DM
>		C desktop.ini	402 B	Fed 1, 2024, 2:48 PM
1				
		4		•

#### National Czech Programme

RStudio File Edit Code View Plots Session Build I	Debug Profile Tools Help		– 0 X
New File	▶ tion 🛛 😳 👻 Addins 🕶		🔋 Project: (None)
New Project		Environment History Connections Tutorial	
New Project         Open File in New Column         Recent Files         Open Project         Open Project in New Session         Recent Projects         Import Dataset         Save       Ctrl+S         Save As         Save All         Alt+Ctrl+S         Print         Close       Ctrl+Shift+W         Close All       Ctrl+Shift+W         Close Project       Quit Session         Quit Session       Ctrl+Q	<pre>%1 ucrt) "Eye Holes" oundation for Statistical Computing v32/x64 (64-bit) nes with ABSOLUTELY NO WARRANTY. ibute it under certain conditions. ice()' for distribution details. sct with many contributors. more information and ie R or R packages in publications. nos, 'help()' for on-line help, or . browser interface to help.</pre>	Environment       History       Connections       Tutorial         Import Dataset •       6 MiB •       Import Dataset         R •       Global Environment •         Environment is emption         Import Dataset •       Import Dataset         Import Dataset •       Import Dataset •         Import Dataset •       Import Dataset •	List - C List - C Modified May 12, 2025, 3:45 PM Feb 1, 2024, 2:48 PM

R	New Directory Start a project in a brand new working directory	>
R	<b>Existing Directory</b> Associate a project with an existing working directory	>
	Version Control Checkout a project from a version control repository	>

• Vyberte složku, do které jste si uložili data, se kterými budeme pracovat.

Back	Create Project from Ex	xisting Directory	
	Project working directory:		Browse.
K			

• Vaše okna by měla po vytvoření projektu vypadat následovně:

National Czech

Programme

Deosc





#### Popis oken

- Rstudio je v základním nastavení rozloženo do čtyř oken.
- Předtím, než si jednotlivá okna popíšeme, je potřeba abyste si vytvořili nový prázdný R skript.
- Skript můžete vytvořit pomocí klávesové zkratky či přímo výběrem, viz. obrázek vpravo.

RStudio						
File Edit Code View Plots	s Session Build Debug Profile Tools Help					
🌔 - 👒 🕣 - 🕞 🔝	📥 🛛 🖈 Go to file/function 👘 🗄 👻 Addins 👻					
R Script Ctrl+Shift+N						
🐑 Quarto Document						
菜 Quarto Presentation	2(2023-10-31) ucrt) "Eve					
R Notebook	2023 The R Foundation for S					
R Markdown	64-w64-mingw32/x64 (64-bit)					
📧 Shiny Web App						
Plumber API	ware and comes with ABSOLUT					
Text File	e to redistribute it under )' or 'licence()' for distr					
💬 C++ File	) or licence() for discr					
Python Script	rative project with many co					
SQL Script	tors()' for more informatio					
🐑 Stan File	n how to cite R or R packa					
D3 Script						
R Sweave	for some demos, 'help()' fo					
횐 R HTML	tor an HTML browser interf					
R Documentation	quit R.					





- Skript
  - zde píšete kód
- Konzole
  - zde vidíte výsledky kódu a programové oznámení (např. Error)
- Globální prostředí
  - seznam načtených dat a přiřazených objektů
- Multifunkční okno
  - Files: seznam souborů v directory ve kterém pracujete
  - Plots: vytvořené grafy
  - Packages: balíčky (pluginy)
  - Help: nápověda



#### Klávesové zkratky

- K psaní kódu budete potřebovat umět napsat několik specifických symbolů.
- Zde můžete najít jejich seznam společně s klávesovými zkratkami pro Windows s českou klávesnicí.
- Doporučuji si psaní symbolů vyzkoušet a během školení mít tento seznam po ruce.
- Pokud používáte jinou klávesnici a symboly proto mají jiné zkratky, doporučuji vám vytvořit si vlastní verzi seznamu zkratek.

National Czech Programme

#### Klávesové zkratky

Symbol	Zkratka	Symbol	Zkratka
#	Alt + x	>	Alt+.
\$	Alt + ů	<	Alt+,
&	Alt + c	{ }	Alt + b/n
~	Alt + +	[]	Alt + f/g
	Alt + w	٨	Alt + 94

# R pro analýzu vědeckých dat

Lucie Hošková

### Odkaz na prezentaci

https://tinyurl.com/EOSC-RproDS

### Organizace

- bloky 4 x 45 minut
  - 1. Úvod práce s R a Rstudiem
  - 2. Nahrávání dat a manipulace s dataframe
  - 3. Statistická analýza a testování hypotéz
  - 4. Vizualizace
- kombinace code-along v RStudiu a prezentace
- mock-up reálných situací práce s daty
- doplňující otázky na konci bloku

### 1. blok

- základní typy hodnot v R
- operátory
- spuštění kódu
- komentáře a pravidla psaní kódu
- vytvoření objektu
- funkce, balíčky
- typy objektů

## Typy objektů

- Skalár
  - objekt obsahující pouze a právě jednu hodnotu

```
1 sk1 <- 5; sk1
[1] 5
1 sk2 <- "ID"; sk2</pre>
```

[1] "ID"

- Vektor
  - objekt složený z x skalárů stejného typu hodnot
  - jednorozměrný objekt

Funkce c() vytvoří vektor z vložených hodnot".

1 cisla <- c(5, 8, 11); cisla

[1] 5 8 11

1 popisky <- c("ID", "No.", "Age"); popisky</pre>

```
[1] "ID" "No." "Age"
```

# Více typů hodnot ve vektoru vede k nucené konverzi na character

		1 mix <- c("ID", 2, TRUE); mix
[1]	"ID"	"2" "TRUE"
		1 is.character(mix)
[1]	TRUE	
		1 c(sk1, sk2) #muzeme kombinovat i objekty
[1]	"5"	"ID"

#### • Data frame

#### dvojrozměrný objekt složený z vektorů stejné délky

		1 df	<- dat	ta.frame(cisla, popisky, mix); df
	cisla	popisky	mix	
1	5	ID	ID	
2	8	No.	2	
3	11	Age	TRUE	

#### Jaká je struktura dataframe? Z jakých vektorů se skládá?

#### 1 str(df)

'data.frame': 3 obs. of 3 variables: \$ cisla : num 5 8 11 \$ popisky: chr "ID" "No." "Age" \$ mix : chr "ID" "2" "TRUE"

## CVIČENÍ

- 1. Překrývají se nějaká čísla v sekvencích 8:20 a 5:12?
- 2. Vytvořte alespoň dvěmi způsoby číselnou řadu 1122334.

## ŘEŠENÍ

1. Překrývají se nějaká čísla v sekvencích 8:20 a 5:12?

1 5:12 %in% 8:20

[1] FALSE FALSE FALSE TRUE TRUE TRUE TRUE TRUE

#### 2. Vytvořte alespoň dvěmi způsoby číselnou řadu 1122334.

1 rep(1:4, each = 2, length.out = 7)

[1] 1 1 2 2 3 3 4

1 c(1,1,2,2,3,3,4)

[1] 1 1 2 2 3 3 4

1 c(rep(1:3, each = 2), 4)

[1] 1 1 2 2 3 3 4

### 2. blok

- načtení dat
- práce s dataframe
- faktory
- výběr dat
- balíček dplyr
- úprava dataframe

### Načtení dat

Načtení tabulkového datasetu do R se provádí pomocí **read.table** funkcí.

- read.delim() pro tabulátorem oddělené hodnoty
- read.csv() pro CSV soubory

Načteme si zkušební dataset newborns.txt.

```
1 data <- read.delim("newborns.txt", header = TRUE, sep = "\t", dec
```

Nešlo by to jednodušeji?

1 data2 <- read.delim("newborns.txt")</pre>

Načetli se datasety stejně? Proč?

1 identical(data, data2)

### Popis a struktura datasetu

#### 1 str(data)

'data.frame': 1402 obs. of 4 variables: \$ edu.M : int 2 2 2 1 3 2 1 2 2 1 ... \$ prch.N : int 0 0 0 0 0 1 2 0 0 0 ... \$ sex.C : chr "m" "m" "f" "m" ... \$ weight.C: int 3470 3240 2980 3280 3030 3650 4080 3040 3070 3110 ...

- edu.M: vzdělání matky
  - 1: základní
  - 2: středoškolské bez maturity
  - 3: středoškolské s maturitou
  - 4: vysokoškolské
- prch.N: počet biologických starších sourozenců
- sex.C: pohlaví dítěte
  - m: muž
  - f: žena
- weight.C: porodní hmotnost dítěte v gramech

## Faktory

# Kategoriální data lze v R zapsat ve specifickém data typu - **faktoru**.

```
1 data$sex.C <- as.factor(data$sex.C)</pre>
```

```
2 levels(data$sex.C)
```

[1] "f" "m"

1 str(data)

```
'data.frame': 1402 obs. of 4 variables:
$ edu.M : int 2 2 2 1 3 2 1 2 2 1 ...
$ prch.N : int 0 0 0 0 1 2 0 0 0 ...
$ sex.C : Factor w/ 2 levels "f", "m": 2 2 1 2 2 2 1 1 1 2 ...
$ weight.C: int 3470 3240 2980 3280 3030 3650 4080 3040 3070 3110 ...
```

# R umožňuje konverzi dat na faktory a zároveň **přejmenování úrovní**.

	<pre>1 data\$edu.M &lt;- factor(data\$edu.M, labels = c("ZS", "SS", "SSm", "VS" 2 levels(data\$edu.M)</pre>
[1] "ZS"	"SS" "SSm" "VS"
	1 head(data)

	edu.M	prch.N	sex.C	weight.C
1	SS	0	m	3470
2	SS	0	m	3240
3	SS	0	f	2980
4	ZS	0	m	3280
5	SSm	0	m	3030
6	SS	1	m	3650

### Výběr dat z objektů pomocí base R

V base R jsou dva hlavní způsoby výběru dat.

- výběr pomocí []
  - objekt[určení výběru]

```
1 vektor <- 2:30; vektor
[1] 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24
25 26
[26] 27 28 29 30
1 vektor[3]</pre>
```

[1] 4

Při použití [] je potřeba specifikovat jakou část chcete vybrat u **každé dimenze**. Dataframe je například dvojrozměrný objekt (má řádky a sloupce), je proto potřeba definovat výběr u obojího.
		1 data	[1, 4]	#1. radek, 4. sloupec							
[1]	3470										
	1 data[c(1, 5, 7),4]										
[1]	[1] 3470 3030 4080										
		1 data	[1:3, ]	#1 3. radek, vsechny sloupce							
e	edu.M prch.N sex.C weight.C										
1	SS	0	m	3470							
2	SS	0	m	3240							
3	SS	0	f	2980							
		1 data	[2, -1]	#2. radek, vsechny sloupce krome 1.							
p	rch.N	sex.C we	eight.C								
2	0	m	3240								

## Lze použít i **jména sloupců**.

		1	data[	1:3, c("e	edu.M", '	"weight	.C")]					
	edu.M	wei	ght.C									
1	SS		3470									
2	SS		3240									
3	SS		2980									
		1	data[	"edu.M"]	#subset	pouze	s jedno	u dimenzi	u dat	taframe	vzdy	vra
	edu	л.М										
1		SS										
2		SS										
3		SS										

4	ZS
5	SSm
6	SS
7	ZS
8	SS
9	SS
10	ZS
11	SSm
12	SSm
13	ZS
14	ZS

## • výběr pomocí \$

## objekt\$název sloupce

	1	data	\$edu.	М										
[1]	SS	SS	SS	ZS	SSm	SS	ZS	SS	SS	ZS	SSm	SSm	ZS	ZS
[15]	SS	ZS	ZS	ZS	SS	SS	SSm	ZS	ZS	SSm	SSm	ZS	ZS	ZS
[29]	SS	SSm	ZS	SS	SSm	ZS	ZS	SSm	ZS	ZS	ZS	ZS	SSm	SS
[43]	ZS	ZS	ZS	ZS	SS	SS	SS	ZS	SS	ZS	SS	SS	SS	ZS
[57]	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	SS	ZS	ZS	ZS
[71]	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	SS	ZS	SS	SSm	SS	ZS
[85]	ZS	ZS	ZS	ZS	SS	SS	ZS	ZS	SS	ZS	ZS	SS	ZS	SS
[99]	ZS	SS	ZS	ZS	SSm	SS	ZS	SS	SS	ZS	ZS	ZS	SS	SS
[113]	ZS	ZS	SS	ZS	ZS	SS	ZS	SS	ZS	ZS	SS	ZS	ZS	ZS
[127]	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS
[141]	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	SS
[155]	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS	ZS

[169]	ZS	ΖS	ZS	ZS	ΖS	ZS	ZS	ΖS	ΖS	ZS	ZS	ZS	SSm	ΖS
[183]	ZS	ZS	SS	ZS	ZS	ZS	SS	ΖS	ΖS	SS	ZS	ZS	ZS	ΖS
[197]	ZS	ZS	SSm	SS	ZS	ZS	ZS	SSm	ZS	ZS	ZS	SS	SS	ΖS

## Rozdíl mezi výsledky výběru

- \$ vrátí *vektor* obsahující hodnoty vybraného sloupce
  - sloupec musí mít jasně určené jméno
  - nelze vybírat řádky
- [] vrátí část *objektu* obsahující vybrané hodnoty
  - zachovává si přitom vlastnosti objektu

```
1 identical(data["edu.M"], data$edu.M)
[1] FALSE
1 mode(data$edu.M)
[1] "numeric"
1 mode(data["edu.M"])
```

[1] "list"

#### Tyto znalosti pak lze použít k výběru **specifických hodnot**.

	1	data[d	ata\$ed	u.M == "ZS",	]
	edu.M	prch.N	sex.C	weight.C	
4	ZS	0	m	3280	
7	ZS	2	f	4080	
10	ZS	0	m	3110	
13	ZS	1	m	2940	
14	ZS	0	f	1550	
16	ZS	1	f	1960	
17	ZS	0	f	2910	
18	ZS	0	f	3670	
22	ZS	1	f	2660	
23	ZS	2	f	2550	
26	ZS	1	f	2540	
27	ZS	0	m	3170	
28	ZS	0	m	2970	
31	ZS	1	f	2160	
<u> </u>		1		2110	
	1	data[d	ata\$ed	u.M == "ZS"	<pre>&amp; data\$prch.N &gt;= 2, c("sex.C", "weight.C"</pre>
	sex.C	weight.(			
7	f	4080	)		
23	f	255(	)		
35	m	2480	)		
87	f	314(	)		
91	f	3040	)		
182	m	2710	)		
186	m	2780	)		
187	m	2190	)		

205	m	3670
207	m	2390
210	f	2540
220	f	2850
235	f	3430
239	f	1590

# Výběr dat z objektů pomocí balíčku {dplyr}

Balíček dplyr je součástí balíčku **tidyverse**, balíček dplyr jsme si nainstalovali během prvního bloku.

1 library(dplyr)

Umožňuje výběr pomocí jednodušších příkazů než těch, které požaduje base R.

	1	select(	data,	edu.M,	sex.C)	)			
	edu.M se	ex.C							
1	SS	m							
2	SS	m							
3	SS	f							
4	ZS	m							
5	SSm	m							
6	SS	m							
7	ZS	f							
8	SS	f							
9	SS	f							

10	ZS	m
11	SSm	f
12	SSm	m
13	ZS	m
14	7.5	f

### Dokáže také provést výběr dle specifických podmínek.

		1 filt	er(data,	edu.M ==	"SS", weight.C > 1500)
	edu.M	prch.N	sex.C we	eight.C	
1	SS	0	m	3470	
2	SS	0	m	3240	
3	SS	0	f	2980	
4	SS	1	m	3650	
5	SS	0	f	3040	
6	SS	0	f	3070	
7	SS	0	m	2040	
8	SS	0	m	2790	
9	SS	0	m	3490	
10	SS	1	m	2460	
11	SS	1	m	2350	
12	SS	0	f	3340	
13	SS	0	f	4060	
14	SS	0	m	3770	
1 🗆	0.0	$\cap$	٢		

### Co všechno balíček zvládne?

1 help(package = dplyr)

# Úprava datasetu

K dataframe je také možno údaje přidávat. Vytvořme se sloupec weight.K, který bude obsahovat zda se porodní váha novorozence nacházela v normě, či zda byla vysoká či nízká.

- nízká: weight. C < 2500
- norma: 2500 <= weight.C <= 4200
- vysoká: weight.C > 4200

```
1 data$weight.K[data$weight.C < 2500] <- "nizka"
2 data$weight.K[data$weight.C >= 2500 & data$weight.C <= 4200] <- "no
3 data$weight.K[data$weight.C > 4200] <- "vysoka"
4
5 head(data)
```

	eau.M	prcu.N	sex.C	weight.c	weight.r
1	SS	0	m	3470	norma
2	SS	0	m	3240	norma
3	SS	0	f	2980	norma
4	ZS	0	m	3280	norma

5	SSm	0	m	3030	norma
6	SS	1	m	3650	norma

#### Hodnoty si převedeme na faktor.

	1 data\$weight.K 2	<- as.fact	tor(data\$weight	.K)
	3 summary(data)	#shrn mi d	objekt data	
edu.M	prch.N	sex.C	weight.C	weight.K
ZS :420	Min. :0.0000	f:673	Min. : 580	nizka : 272
SS :451	1st Qu.:0.0000	m:729	1st Qu.:2670	norma :1080
SSm :437	Median :1.0000		Median :3170	vysoka: 44
VS : 81	Mean :0.9492		Mean :3071	NA's : 6
NA's: 13	3rd Qu.:1.0000		3rd Qu.:3560	
	Max. :9.0000		Max. :4970	
	NA's :5		NA's :6	

V summary vidíme, že dataframe má v sobě **jedince s chybějícími hodnotami**. Vytvoříme si proto soubor čistě s úplnými údaji.

	1 data <- na.omi 2 summary(data)	t(data)		
edu.M	prch.N	sex.C	weight.C	weight.K
ZS :417	Min. :0.0000	f:662	Min. : 580	nizka : 266
SS :448	1st Qu.:0.0000	m:719	1st Qu.:2680	norma :1071
SSm:435	Median :1.0000		Median :3170	vysoka: 44

VS	:	81	Mean	:0.9428	Mean	:3078
			3rd Qu.	:1.0000	3rd Qu.	:3570
			Max.	:9.0000	Max.	:4970

# CVIČENÍ

Vyberte data pro všechny chlapce, kteří se narodili s porodní váhou mezi 2000 a 2200 gramy (včetně) matkám prvorodičkám s nižším stupněm vzdělání než je vysokoškolský titul.

# ŘEŠENÍ

### První možnost

	1 2 3 4	data[	data\$w data data data data	eight.C % \$sex.C == \$edu.M != \$prch.N =	in% 2000:2 "m" & #ch "VS" & #v = 0,] #0 s	200 & #vaha v intervalu 2000:2200 lapci zdelani krome VS tarsich sourozencu
	edu.M	prch.N	sex.C	weight.C	weight.K	
15	SS	0	m	2040	nizka	
290	SS	0	m	2120	nizka	
299	SSm	0	m	2180	nizka	
398	SS	0	m	2010	nizka	
403	SSm	0	m	2190	nizka	
832	SSm	0	m	2150	nizka	
842	ZS	0	m	2160	nizka	
850	ZS	0	m	2020	nizka	
961	SSm	0	m	2090	nizka	
1290	SSm	0	m	2170	nizka	

### Druhá možnost

	4 5			data\$edu data\$prc	.M %in% c("ZS", ch.N == 0,]	"SS",	"SSm")	& #vycet	vsech	VZC
	edu.M	prch.N	sex.C	weight.C	weight.K					]
15	SS	0	m	2040	nizka					
290	SS	0	m	2120	nizka					
299	SSm	0	m	2180	nizka					
398	SS	0	m	2010	nizka					
403	SSm	0	m	2190	nizka					
832	SSm	0	m	2150	nizka					
842	ZS	0	m	2160	nizka					
850	ZS	0	m	2020	nizka					
961	SSm	0	m	2090	nizka					
1290	SSm	0	m	2170	nizka					

### Třetí možnost

		1 fil 2 3 4	ter(data, sex.( edu.N prch.	weight C == "m" A != "VS N == 0)	.C %in% 20 , ",	00:2200,		
	edu.M prch.N sex.C weight.C weight.K							
1	SS	0	m	2040	nizka			
2	SS	0	m	2120	nizka			
3	SSm	0	m	2180	nizka			
4	SS	0	m	2010	nizka			
5	SSm	0	m	2190	nizka			
6	SSm	0	m	2150	nizka			
7	ZS	0	m	2160	nizka			
8	ZS	0	m	2020	nizka			

9	SSm	0	m	2090	nizka
10	SSm	0	m	2170	nizka

# Jaké jsou mezi výsledky rozdíly? Kdy je lepší použít který přístup?

## 3. blok

- popisná statistika
- dplyr::summarise
- dplyr::group\_by()
- pipes
- analýza nominálních dat
- analýza ordinálních dat
- analýza intervalových dat

## Popisná statistika

## Jakými různými způsoby můžeme získat základní statistické údaje o objektu? Známe již funkci **summary()**.

	1 summary(data)			
edu.M	prch.N	sex.C	weight.C	weight.K
ZS :417	Min. :0.0000	f:662	Min. : 580	nizka : 266
SS :448	1st Qu.:0.0000	m:719	1st Qu.:2680	norma :1071
SSm:435	Median :1.0000		Median :3170	vysoka: 44
VS : 81	Mean :0.9428		Mean :3078	
	3rd Qu.:1.0000		3rd Qu.:3570	
	Max. :9.0000		Max. :4970	

# Jak můžeme tyto údaje získat **individuálně**? Ukažme si to pomocí sloupce weight.C.

		1	<pre>mean(data\$weight.C) #prumer</pre>
[1]	3078	.027	7
		1	<pre>min(data\$weight.C) #minimum</pre>
[1]	580		
		1	max(data\$weight.C) #maximum

```
[1] 4970
```

	1 median(data\$weight.C) #median
[1] 3170	
	1 quantile(data\$weight.C, type = 2) #kvartily
0% 25% 580 2680	50% 75% 100% 3170 3570 4970
	1 quantile(data\$weight.C, probs = 0.25, type = 2) #pouze prvni kvarti
25% 2680	
	1 IQR(data\$weight.C, type = 2) #interkvartilove rozpeti
[1] 890	

## Funkce summary nám ukáže i zastoupení jednotlivých faktorů u slupců s kategorickými hodnotami. Stejný výsledek dostaneme i pomocí funkce **table()**.

```
1 table(data$edu.M)
ZS SS SSm VS
417 448 435 81
1 table(data$edu.M, data$sex.C)
```

ZS 190 227 SS 212 236 SSm 212 223 VS 48 33

Tato funkce vytvoří *kontingenční tabulku* absolutních četností faktorů vybraných hodnot.

Výpočet počtu zastoupení můžeme provést i jednotlivě, a to funkcí **sum()**.

1 sum(data\$edu.M == "ZS")

[1] 417

Funkce sum() spočítá *kolikrát se objeví* x v daném datovém objektu.

Podobnou tabulku jako dostaneme z table() si poté můžeme vytvořit pomocí funkce **data.frame()**.

```
1 tabulka_cetnosti_edu.M <- data.frame(sum(data$edu.M == "ZS"),
2 sum(data$edu.M == "SS"),
3 sum(data$edu.M == "SSm"),
4 sum(data$edu.M == "VS")); tabulka cetnosti edu.M
```

```
sum.data.edu.M....ZS.. sum.data.edu.M....SS.. sum.data.edu.M....SSm..

1 417 448 435

sum.data.edu.M...VS..

1 81
```

Tato funkce vytvoří *dataframe*, kdy po sloupcích spojí vložené objekty. Pokud bychom chtěli spojit objekty po řádcích, můžeme použít funkci **rbind.data.frame()**.

Vidíme, že R **automaticky vygenerovalo názvy sloupců** na základě jejich vlastností. Pro lepší přehled si proto sloupce přejmenujeme, a to pomocí funkce **colnames()**.

1 colnames(tabulka\_cetnosti\_edu.M) <- c("ZS", "SS", "SSm", "VS")</pre>

Podobnou funkci má u řádků funkce row.names().

Této potřebě přejmenovávání se dá vyhnout tak, že do funkce data.frame() rovnou vložíme **pojmenované objekty**.

```
1 ZS <- sum(data$edu.M == "ZS")
2 SS <- sum(data$edu.M == "SS")
3 SSm <- sum(data$edu.M == "SSm")
4 VS <- sum(data$edu.M == "VS")</pre>
```

5 6 data.frame(ZS, SS, SSm, VS)

ZS SS SSm VS

1 417 448 435 81

# dplyr::summarise

K vytváření popisných tabulek lze také použít funkci z balíčku dplyr, **summarise()**. Pro použití funkce musíte mít buď aktivovanou knihovnu balíčku (jako již máme), nebo funkci můžete zavolat přímo bez aktivace knihovny pomocí **dplyr::summarise()**.

```
1 dplyr::summarise(data, prumer_vahy = mean(weight.C),
2 max_sourozencu = max(prch.N))
```

```
prumer_vahy max_sourozencu
1 3078.027 9
```

Funkce umožňuje vytvářet dataframe pomocí pojmenování statistických funkcí aplikovaných na hodnoty definovaného objektu.

Pro více informací o tom, co vše funkce dokáže, si můžete

# dplyr::group\_by()

Při práci s funkcí summarise se můžete setkat se situací, kdy byste rádi rozdělili výsledný dataframe podle určitých kategorií. K tomu slouží funkce **group\_by()**.

1 m\_vzdelani <- group\_by(data, edu.M)</pre>

**Funkce sama objekt nezmění**, ale při kombinaci s funkcí summarise umožní rozdělení statistických analýz dle kategorií definovaných v group\_by.

	1	summarise(	(m_vzdelani, prume)	_vahy = mean(weight.C),		
	2		max_souroze	encu = max(prch.N))		
#	# A tibble: 4 × 3					
	edu.M pru	mer_vahy ma	x_sourozencu			
	<fct></fct>	<dbl></dbl>	<int></int>			
1	ZS	2931.	9			
2	SS	3128.	6			
3	SSm	3153.	4			
4	VS	3157.	3			

## Pipes

Je možné, že budete provádět analýzu na výběru z výběru z výběru datasetu, jako se stalo v přechozích kapitolách. Jak si práci usnadníte, abyste nemuseli stále ukládat další a další objekty? Použitím takzvaného **operátoru pipe** %>%. Jeho funkce bud ukázána na následujícím kódu, jednou napsaném bez použití operátoru a jednou s ním:

1	matka_ZS <- filter(data, edu.M == "ZS")
2	vaha_srovnana <- arrange(matka_ZS, weight.C
3	head(vaha srovnana, <mark>5</mark> )

edu.M prch.N sex.C weight.C weight.K

	T			2	2
1	ZS	0	f	750	nizka
2	ZS	2	m	970	nizka
3	ZS	0	f	970	nizka
4	ZS	0	m	1190	nizka
5	ZS	0	m	1210	nizka

V tomto příkladě se setkáváme s novou funkcí **dplyr::arrange()**, která seřadí data v X podle údajů v Xa od nejnižší hodnoty po nejvyšší.

		1 da	ta %>%		
		2	filter(e	edu.M ==	"ZS") %>%
		3	arrange	(weight.	C) 응>응
		4	head( <mark>5</mark> )		
	edu.M	prch.N	sex.C w	eight.C	weight.K
1	ZS	0	f	750	nizka
2	ZS	2	m	970	nizka
3	ZS	0	f	970	nizka
4	ZS	0	m	1190	nizka
5	ZS	0	m	1210	nizka

# CVIČENÍ

Spočítejte, jaký je medián a IQR váhy novorozenců, jejichž matky mají už alespoň jednoho potomka a dosáhly nanejvýše středního vzdělání s maturitou. Berte ohled na pohlaví. Výsledky uložte do tabulky.

# Řešení - první možnost

```
1 chlapci <- filter(data, prch.N >= 1,
                         edu.M != "VS", sex.C == "m")
      2
      3
      4 divky <- filter(data, prch.N >= 1,
      5
                          edu.M != "VS", sex.C == "f")
      6
     7
        median chlapci <- median(chlapci$weight.C)</pre>
        IQR chlapci <- IQR(chlapci$weight.C, type = 2)</pre>
      8
      9
        median divky <- median(divky$weight.C)</pre>
    10
        IQR divky <- IQR(divky$weight.C, type = 2)</pre>
    11
    12
    13
        data.frame (median chlapci,
    14
                          IQR chlapci,
                          median divky,
    15
    16
                          IQR divky)
1 4
```

	median_	_chlapci	IQR_chlapci	median_divky	IQR_divky
1		3290	860	3155	830

## Řešení - druhá možnost



# A tibble: 2 × 3
 sex.C median IQR
 <fct> <dbl> <dbl>
1 f 3155 830
2 m 3290 860

# Analýza dat

Poté co jste si upravili dataset a seznámili jste se s ním, tak na něm můžete provádět různé analýzy.

## Nominální data

- kvalitativní data
- barva očí, typ publikace
- nelze je hierarchicky porovnávat

Ordinální

- kvalitativní data
- pořadí narození jedinců, známky ve škole
- nelze určit míru rozdílu mezi stupni

Intervalová

- kvantitativní data
- po sobě jdoucí data s jasnými intervaly
- výška, počet dětí

# Analýza nominálních dat

## Cramerův koeficient

<0:1>

- 0: nízká závislost
- 1: vysoká závislost

1 porovnani <- table(data\$sex.C, data\$edu.M)</pre>

2 lsr::cramersV(porovnani) #crameruv koeficient

[1] 0.06183871

1 chisq.test(porovnani) #ziskani p-value

Pearson's Chi-squared test

```
data: porovnani
X-squared = 5.281, df = 3, p-value = 0.1523
```

## Pozor! Funkce cramersV je určená specificky k **porovnávání kontingenčních tabulek**!

Mezi vzděláním matky pohlavím dítěte je pouze velmi nízká statisticky nevýznamná závislost.

# Analýza ordinálních dat

Spearmanův korelační koeficient

<-1:1>

- -1: nízká přímá závislost
- 1: vysoká přímá závislost

Lze jej použít i při **porovnávání ordinálních a intervalových dat**, protože je postavený na provnávání hodnot s jasnými rozestupy - například první stupeň je přesně o jednu hodnotu níže než druhý a 99.8 je přesně o dvě hodnoty níže než 100.

1 cor.test(data\$prch.N, data\$weight.C, method = "spearman")

Spearman's rank correlation rho

data: data\$prch.N and data\$weight.C

Mezi počtem starších sourozenců a porodní hmotností je jen velmi nízká přímá závislost.

Vzhledem k p-value vyšší jak 0.05 můžeme prohlásit, že velmi nízká závislost mezi počtem starších sourozenců a porodní hmotností je statisticky nevýznamná.

# Analýza intervalových dat

Pearsonův korelační koeficient

<0:1>

- 0: nízká přímá závislost
- 1: vysoká přímá závislost

Jelikož náš dataset neobsahuje dvoje intervalová data k porovnání, použijeme testovací dataset trees který je v balíčku base R.

		1 he	ad(trees)
	Girth	Height	Volume
1	8.3	70	10.3
2	8.6	65	10.3
3	8.8	63	10.2
4	10.5	72	16.4

5 10.7 81 18.8

6 10.8 83 19.7

1 cor.test(trees\$Girth, trees\$Height, method = "pearson")

Pearson's product-moment correlation

```
data: trees$Girth and trees$Height
t = 3.2722, df = 29, p-value = 0.002758
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
    0.2021327 0.7378538
sample estimates:
        cor
0.5192801
```

## Lze použít i bez určení argumentu method, jelikož Pearson se bere jako baseline.

1 cor.test(trees\$Girth, trees\$Height)

```
Pearson's product-moment correlation
```

```
data: trees$Girth and trees$Height
t = 3.2722, df = 29, p-value = 0.002758
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
   0.2021327 0.7378538
sample estimates:
```
cor 0.5192801

Mezi obvodem a výškou kmene stromů v datasetu je středně pozitivní závislost.

Vzhledem k tomu že p-value < 0.05 můžeme říct že závislost mezi hodnotami je statisticky významná.

### 4. blok

- úvod do vizualizace v R
- balíček ggplot2
- barplot
- histogram
- boxplot
- scatterplot

### Vizualizace v R

# R je schopné provádět rozličné typy datových vizualizací, od jendoduchých histogramů

#### 1 library(ggplot2)

```
2
 3 ggplot(mtcars, aes(x = mpg, fill = factor(cyl))) +
     geom histogram(binwidth = 2, color = "black", alpha = 0.7) +
 4
     scale fill manual(values = c("4" = "skyblue", "6" = "orange", "
 5
     labs(title = "Histogram of Miles Per Gallon (mpg)",
 6
          x = "Miles Per Gallon (mpg)",
 7
 8
          y = "Frequency",
 9
          fill = "Number of Cylinders") +
10
     theme minimal() +
     theme(plot.title = element text(hjust = 0.5, size = 20, face =
11
12
           axis.title = element text(size = 14),
           legend.position = "top")
13
```







#### po 3D grafy.



## {ggplot2}

Zápis kódu v ggplot2 je rozděleno do tří základních částí.

- 1. data: definice dat které jsou vizualizovány
- 2. (aes): mapování mezi hodnotami v datsetu a vizuální styl
- 3. geom(): část definující jaký typ vizualizace má být vytvořena.

Tyto části jsou ty **základní, které jsou k fungování ggplot nutné**, ale **je možné k nim přidat** nepřeberné množství dalších podfunkcí, které tak mohou dále vizualizaci vyšperkovat.

### Barplot

### Určený k porovnávání jednotlivých kategorií a jejich variant.

```
1 library(ggplot2)
2 ggplot(data = data) +
3 aes(x = edu.M) + #y neni definovano tim padem R secte jednotlive
4 geom bar() #vykresli jako barplot
```



#### ggplot umožňuje i další vizuální upřesnění grafu.

```
1 ggplot(data = data) +
2 aes(x = edu.M) +
3 geom_bar(width = 0.5, #sire sloupcu
4 color="darkblue", #obrys sloupcu
5 fill=rgb(0.1,0.8,0.1,0.7)) + #vypln sloupcu
6 #rgb funkce k urceni specificke barvy
```



#### Barplot je možné taky rozdělit podle jednotlivých skupin.



Počty matek dle vzdělání a pohlaví novorozenců



## Histogram

### Určený k vizualizaci rozložení četnosti veličiny.



# R automaticky vybere počet sloupců. Ten si ovšem můžeme upravit, a to pomocí **Sturgessova pravidla**.

1 round(1 + 3.3 \* log10(length(data\$weight.C)))

[1] 11

1 #log10(pocet hodnot)

### Vytvoříme si histogram znovu se správným počtem sloupců a rovnou si k němu přidáme **vizualizaci hustoty** a **průměrnou hodnotu.**

```
gqplot(data) +
 1
     aes(x = weight.C) +
 2
 3
     geom histogram(aes(y = after stat(density)), #aes prvni cast ko
                     color = "black", fill = "white", bins = 11) + #b
 4
     geom density(alpha=.2, #prusvitnst
 5
                   fill="#FF66666") + #barva vyplne
 6
     geom vline(aes(xintercept = mean(weight.C)), #urceni prumeru
 7
 8
               color = "blue", #barva linie
               linetype = "dashed", #typ linie
 9
               linewidth = 1) #velikost linie
10
```



Můžeme vytvářet i navzájem **překrývající se histogramy**. Všimněte si, že spousta argumentů se nám *opakuje* napříč grafy (například, color, fill, etc).



### Skupiny ale také můžeme rozložit na více grafů.

```
1 ggplot(data,
2 aes(x = weight.C)) +
3 geom_histogram(bins = 11, fill = "pink", color = "red") +
4 facet_grid(sex.C ~ .) #rozdel dle kategorie sex.C do dvou grafu
```



### Boxplot

Určený k vizualizaci rozložení kvantitativních dat. Specifikuje medián, IQR, minimum, maximum a odchylky nad intervalem spolehlivosti.

```
1 ggplot(data) +
2 aes(x = sex.C, #jak bude dle sex.C
3 y = weight.C, #vypadat rozlozeni weight.C?
4 fill = sex.C) +
5 geom_boxplot()
```



# Můžeme si také **vizualizovat jednotlivé hodnoty** přímo na boxplotu.

```
1 ggplot(data) +
2 aes(x = sex.C, #jak bude dle sex.C
3 y = weight.C, #vypadat rozlozeni weight.C?
4 fill = sex.C) +
5 geom_boxplot(notch = TRUE) + #pridame zarez
```





### Scatterplot

### Slouží k porovnávání dvou kvantitativních hodnot.

1 ggplot(data = trees)+

- 2 aes(x = Girth, y = Height) +
- 3 geom\_point(shape = 17)



# Můžeme si upřesnit, aby **velikost bodu** byla závislá na jiné proměnné.





# Taktéž je možné znázornit si graficky **regresní linii** (pokud jsou data lineárně závislá).





### Je možné scatterplot taky vyjádřit pomocí hustoty.

```
1 graf <- ggplot(data = trees) +</pre>
```

```
2 aes(x = Girth, y = Height) +
```

- 3 geom\_point() +
- 4 geom\_density\_2d(); graf



## Uložení grafu

### Graf si můžete ulžit buď skrz okno Plots, a nebo pomocí kódu.

1 ggsave(plot = graf, 2 filename = "ukazka\_grafu.jpg", 3 width = 10, 4 height = 5, 5 units = "cm")

# CVIČENÍ

Vytvořte vizualizaci váhy novorozenců s maximálně jedním starším sourozencem s ohledem na jejich pohlaví. Přejmenujte legendu, popisky os a pojmenujte graf.

## ŘEŠENÍ - první možnost





## ŘEŠENÍ - druhá možnost







### Co dál?

Rozcestník ale úplně na všechno

The Big Book of R

### Knihy (online)

**R for Data Science** 

**R** Graphics Cookbook

Hands-on Programming with R

### Balíčky

tidyverse ggplot2

#### Cheatsheets

base R, base R no.2 dplyr



# Děkuji za pozornost.

lucie.hoskova@ruk.cuni.cz









VŠB TECHNICKÁ | IT4INNOVATIONS |||| UNIVERZITA | NÁRODNÍ SUPERPOČÍTAČOVÉ OSTRAVA | CENTRUM

**Registrační číslo IPs EOSC-CZ** CZ.02.01.01/00/22\_004/0007682