# Polish Language(s) and Digital Humanities Using R

G. Moroz

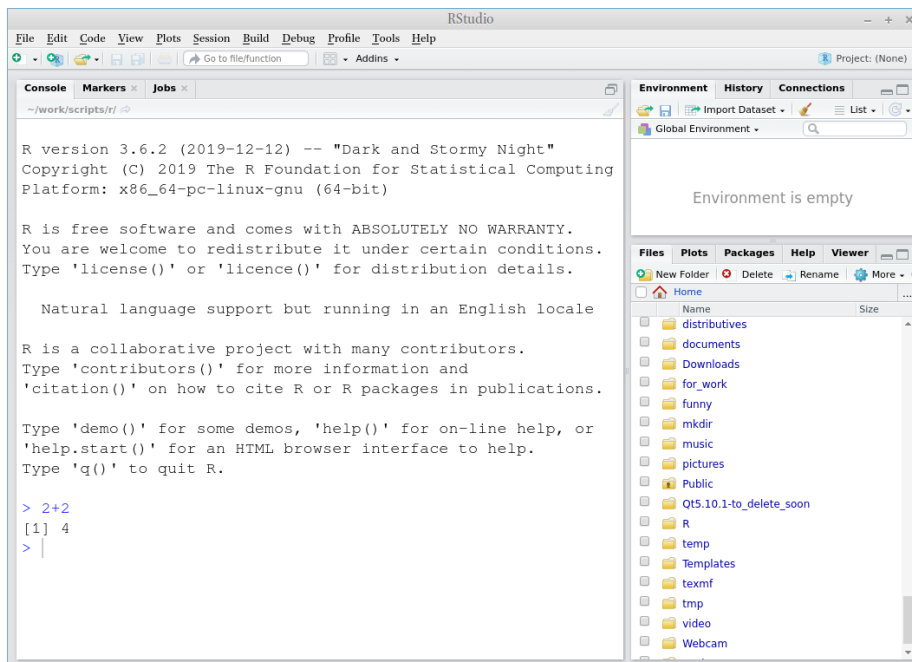2020

ii

# Contents

# Chapter 1

# Prerequisites

Before the classes I would like to ask you to follow the instructions mentioned below to prepare your device for the class work:

- install **R** from the following link: https://cloud.r-project.org/
- install **RStudio** from the following link: https://rstudio.com/products/rstudio/download/#download (FREE version, no need to pay!)
- after the installation run the RStudio program, type `2+2`, and press `Enter`.



If you see something like this, then you are well prepared for classes.

- Go to the https://rstudio.cloud/ website and sign up there. This is optional, but it will be a backup version, if something will not work on your computer.

# Chapter 2

# Introduction to R and RStudio

## 2.1 Introduction

### 2.1.1 Why data science?

Data science is a new field that actively developing lately. This field merges computer science, math, statistics, and it is hard to say how much science in data science. In many scientific fields a new data science paradigm arises and even forms a new sub-field:

- Bioinformatics
- Crime data analysis
- Digital humanities
- Data journalism
- Data driven medicine
- ...

There are a lot of new books "Data Science for ...":

- psychologists (**?**)
- immunologists (**?**)
- business (**?**)
- public policy (**?**)
- fraud detection (**?**)
- ...

Data scientist need to be able:

- gather data
- transform data

- visualize data
- create a statistical model based on data
- share and represent the results of this work
- organize the whole workflow in the reproducible way

### 2.1.2  Why R?

R (**?**) is a programming language with a big infrastructure of packages that helps to work in different fields of science and computer technology.

There are several alternatives:

- Python (**??**)
- Julia (**?**)
- bash (**?**)
- java (**?**)
- …

You can find some R answers here:

- R for data science (**?**), it is online
- R community
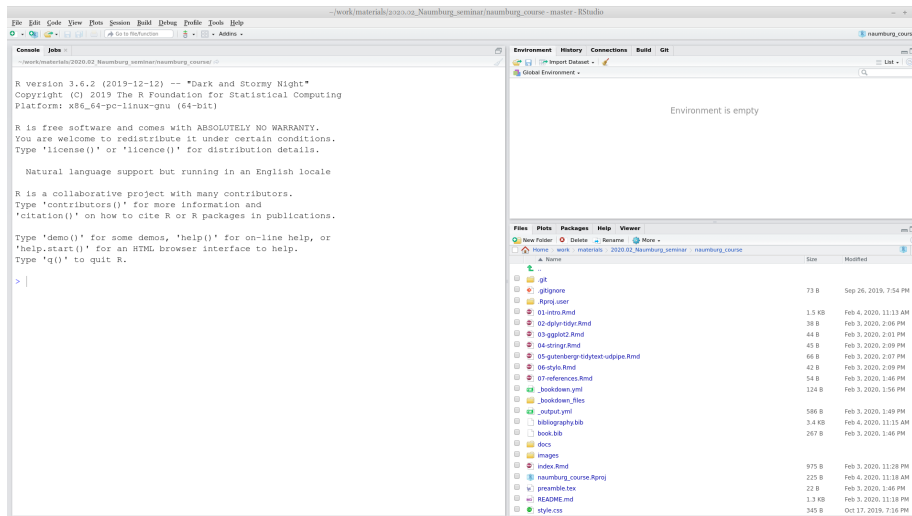- stackoverflow
- any search engine you use
- …

## 2.2  Introduction to RStudio

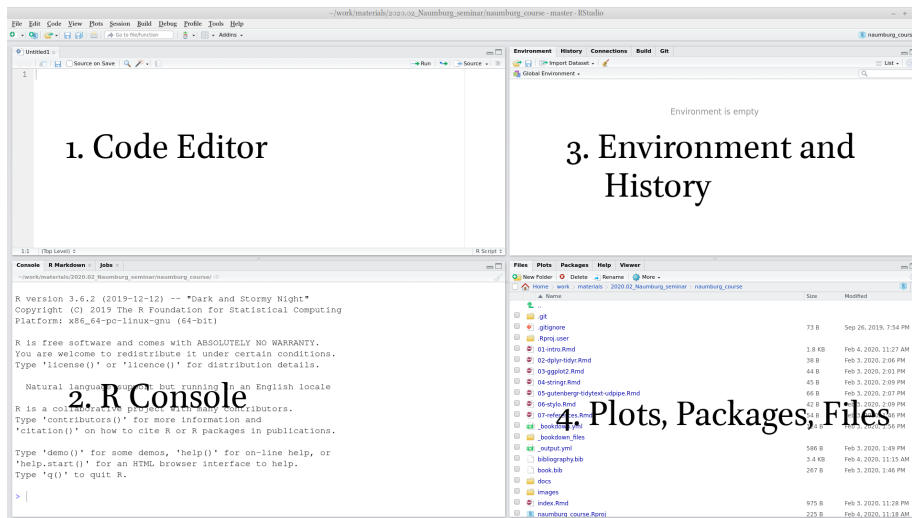R is the programming language. RStudio is the most popular IDE (Integrated Development Environment) for R language.

When you open RStudio for the first time you can see something like this:

When you press  button at the top of the left window you will be able to see all four panels of the RStudio.



## 2.3   R as a calculator

Lets first start with a calculator. Press in R console

```
2+9
```

```
## [1] 11
```

```r
50*(9-20)
```

```
## [1] -550
```

```r
3^3
```

```
## [1] 27
```

```r
9^0.5
```

```
## [1] 3
```

```r
9+0.5
```

```
## [1] 9.5
```

```r
9+.5
```

```
## [1] 9.5
```

```r
pi
```

```
## [1] 3.141593
```

Reminder after division

```r
10 %% 3
```

```
## [1] 1
```

So you are ready to solve some really hard equations (round it four decimal places):

$$\frac{\pi + 2}{2^{3-\pi}}$$

list of hints

Are you sure that you rounded the result?  I expect the answer to be rounded to four decimal places: `0.87654321` becomes `0.8765`.

Are you sure you didn't get into the brackets trap?  Even though there is no any brackets in the mathematical notation, you need to add them in R, otherwise the operation order will be wrong.

## 2.4 Comments

All text after the hash **#** within the same line is considered a comment.

```r
2+2 # it is four
```

```
## [1] 4
```

```r
# you can put any comments here
3+3
```

```
## [1] 6
```

## 2.5 Functions

The most important part of R is functions: here are some of them:

```r
sqrt(4)
```

```
## [1] 2
```

```r
abs(-5)
```

```
## [1] 5
```

```r
sin(pi/2)
```

```
## [1] 1
```

```r
cos(pi)
```

```
## [1] -1
```

```r
sum(2, 3, 9)
```

```
## [1] 14
```

```r
prod(5, 3, 9)
```

```
## [1] 135
```

```r
sin(cos(pi))
```

```
## [1] -0.841471
```

Each function has a name and zero or more arguments. All arguments of the function should be listed in parenthesis and separated by comma:

```r
pi
```

```
## [1] 3.141593
```

```r
round(pi, 2)
```

```
## [1] 3.14
```

Each function's argument has its own name and serial number. If you use names of the function's arguments, you can put them in any order. If you do not use names of the function's arguments, you should put them according the serial number.

```r
round(x = pi, digits = 2)
```

```
## [1] 3.14
```

```r
round(digits = 2, x = pi)
```

```
## [1] 3.14
```

```r
round(x = pi, d = 2)
```

```
## [1] 3.14
```

```r
round(d = 2, x = pi)
```

```
## [1] 3.14
```

```r
round(pi, 2)
```

```
## [1] 3.14
```

```r
round(2, pi) # this is not the same as all previouse!
```

```
## [1] 2
```

There are some functions without any arguments, but you still should use paren-
thesis:

```r
Sys.Date() # correct
```

```
## [1] "2020-02-05"
```

```r
Sys.Date # wrong
```

```
## function ()
## as.Date(as.POSIXlt(Sys.time()))
## <bytecode: 0x625bf08cc528>
## <environment: namespace:base>
```

Each function in R is documented. You can read its documentation typing
question mark before the function name:

```r
?Sys.Date
```

Explore the function `log()` and calculate the following logarithm:

$$\log_3(3486784401)$$

list of hints

A-a-a! I don't remember anything about logarithms... The logarithm is the
inverse function to exponentiation. That means the logarithm of a given number
$x$ is the exponent to which another fixed number, the base $b$, must be raised, to
produce that number $x$.
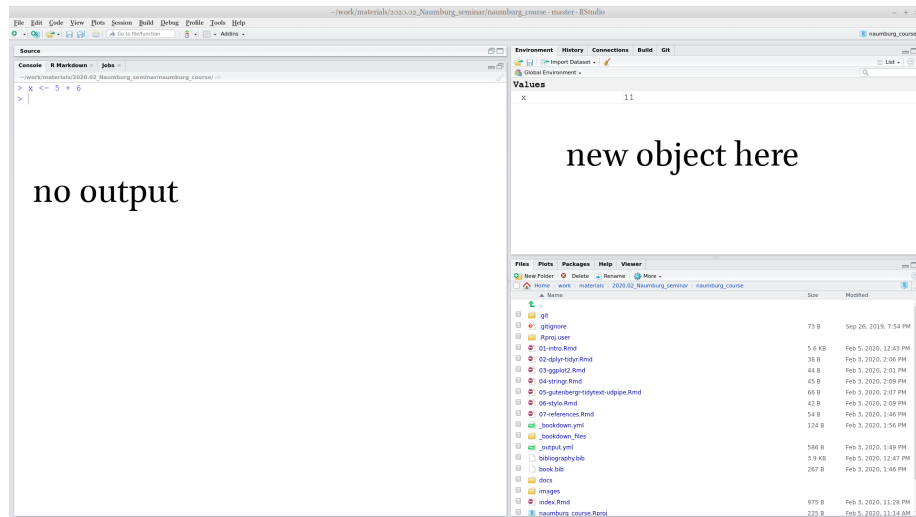
$$10^n = 1000, \text{ what is n?}$$

$$n = \log_{10}(1000)$$

What is this small 3 in the task means? This is the base of the logarithm.
So the task is: what is the exponent to which another fixed number, the base *3*,
must be raised, to produce that number *3486784401*.

## 2.6    Variables

Everything in R can be stored in a variable:

```r
x <- 5 + 6
```



As a result, no output in the Console, and a new variable $x$ appear in the Environment window. From now on I can use this new variable:

```r
x + x
```

```
## [1] 22
```

```r
sum(x, x, 7)
```

```
## [1] 29
```

All those operation don't change the variable value. In order to change the variable value you need to make a new assignment:

```r
x <- 5 + 6 + 7
```

The fast way for creating `<-` in RStudio is to press `Alt -` on your keyboard.

It is possible to use equal sign `=` for assignment operation, but the recommendations are use arrow `<-` for the assignment, and equal sign `=` for giving arguments' value inside the functions.

For removing vector you need to use function `rm()`:

```r
rm(x)
x
```

```
## Error in eval(expr, envir, enclos): object 'x' not found
```

### 2.6.1   Variable comparison

It is possible to compare different variables

```r
x <- 18
x > 18
```

```
## [1] FALSE
```

```r
x >= 18
```

```
## [1] TRUE
```

```r
x < 100
```

```
## [1] TRUE
```

```r
x <= 18
```

```
## [1] TRUE
```

```r
x == 18
```

```
## [1] TRUE
```

```r
x != 18
```

```
## [1] FALSE
```

### 2.6.2   Variable types

There are several types of variable in R. In this course the only important types will be `double` (all numbers), `character` (or strings), and `logical`:

```r
x <- 2+3
typeof(x)
```

```
## [1] "double"
```

```r
y <- "Cześć"
typeof(y)
```

```
## [1] "character"
```

```r
z <- TRUE
typeof(z)
```

```
## [1] "logical"
```

## 2.7   Vector

R object that contain multiple values of the same type is called **vector**. It could be created with the command `c()`:

```r
c(3, 0, pi, 23.4, -53)
```

```
## [1]    3.000000    0.000000    3.141593   23.400000  -53.000000
```

```r
c("Kraków", "Warszawa", "Cieszyn")
```

```
## [1] "Kraków"    "Warszawa" "Cieszyn"
```

```r
c(FALSE, FALSE, TRUE)
```

```
## [1] FALSE FALSE  TRUE
```

```r
a <- c(2, 3, 4)
b <- c(5, 6, 7)
c(a, b)
```

```
## [1] 2 3 4 5 6 7
```

For the number sequences there is an easy way:

```r
1:10
```

```
##  [1]  1  2  3  4  5  6  7  8  9 10
```

```r
3:-5
```

```
## [1]  3  2  1  0 -1 -2 -3 -4 -5
```

From now you can understand that all we have seen before is a vector of length one. That is why there is `[1]` in all outputs: it is just an index of elements in vector. Have a look here:

```
1:60
```

```
##  [1]  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 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50
## [51] 51 52 53 54 55 56 57 58 59 60
```

```
60:1
```

```
##  [1] 60 59 58 57 56 55 54 53 52 51 50 49 48 47 46 45 44 43 42 41 40 39 38 37 36
## [26] 35 34 33 32 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11
## [51] 10  9  8  7  6  5  4  3  2  1
```

There is also a function `sec()` for creation of arithmetic progressions:

```
1:20
```

```
##  [1]  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
```

```
seq(from = 1, to = 20, by = 1)
```

```
##  [1]  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
```

```
seq(from = 2, to = 100, by = 13)
```

```
## [1]   2 15 28 41 54 67 80 93
```

> Use argument `length.out` of function `seq()` and create an arithmetic sequence from $\pi$ to $2\pi$ of length 50.

### 2.7.1   Vector coercion

Vectors are R objects that contain multiple values of **the same type**. But what if we merge together different types?

```
c(1, "34")
```

```
## [1] "1"  "34"
```

```r
c(1, TRUE)
```

```
## [1] 1 1
```

```r
c(TRUE, "34")
```

```
## [1] "TRUE" "34"
```

It is clear that there is hierarchy: strings > double > logical. It is not universal across different programming languages. It doesn't correspond to amount of values of particular type:

```r
c(1, 2, 3, "34")
```

```
## [1] "1"  "2"  "3"  "34"
```

```r
c(1, TRUE, FALSE, FALSE)
```

```
## [1] 1 1 0 0
```

The same story could happen during other operations:

```r
5+TRUE
```

```
## [1] 6
```

### 2.7.2  Vector operations

All operation that we discussed earlier could be done with vectors of the same length:

```r
1:5 + 6:10
```

```
## [1]  7  9 11 13 15
```

```r
1:5 - 6:10
```

```
## [1] -5 -5 -5 -5 -5
```

```r
1:5 * 6:10
```

```
## [1]  6 14 24 36 50
```

There are operation where the vector of any length and vector of length one is involved:

```
1:5 + 7
```

```
## [1]  8  9 10 11 12
```

```
1:5 - 7
```

```
## [1] -6 -5 -4 -3 -2
```

```
1:5 / 7
```

```
## [1] 0.1428571 0.2857143 0.4285714 0.5714286 0.7142857
```

There are a lot of functions in R that are **vectorised**. That means that applying this function to a vector is the same as apply this function to each ellement of the vector:

```
sin(1:5)
```

```
## [1]  0.8414710  0.9092974  0.1411200 -0.7568025 -0.9589243
```

```
sqrt(1:5)
```

```
## [1] 1.000000 1.414214 1.732051 2.000000 2.236068
```

```
abs(-5:3)
```

```
## [1] 5 4 3 2 1 0 1 2 3
```

### 2.7.3 Indexing vectors

How to get some value or banch of values from a vector? You need to index them:

```
x <- c(3, 0, pi, 23.4, -53)
y <- c("Kraków", "Warszawa", "Cieszyn")

x[4]
```

```
## [1] 23.4
```

```r
y[2]
```

```
## [1] "Warszawa"
```

It is possible to have a vector as index:

```r
x[1:2]
```

```
## [1] 3 0
```

```r
y[c(1, 3)]
```

```
## [1] "Kraków"  "Cieszyn"
```

It is possible to index something that you **do not** want to see in the result:

```r
y[-2]
```

```
## [1] "Kraków"  "Cieszyn"
```

```r
x[-c(1, 4)]
```

```
## [1]    0.000000    3.141593  -53.000000
```

### 2.7.4  `NA` value

## 2.8   Dataframe (tibble)

### 2.8.1   Indexing dataframes

## 2.9   Packages

## 2.10   Data import

### 2.10.1   `.csv` files

### 2.10.2   `.xls` and `.xlsx` files

## 2.11   Rmarkdown

# Chapter 3

# Data manipulation: `dplyr`

# Chapter 4

# Data visualisation: `ggplot2`

# Chapter 5

# Strings manipulation: `stringr`

# Chapter 6

# Text manipulation: `gutenbergr`, `tidytext`, `udpipe`

# Chapter 7

# Stylometric analysis: `stylo`