



Data Analysis for the Social Sciences with R

Introduction

Prof. Kevin Koehler
kevin.koehler@santannapisa.it

The class

Aims:

1. Introduce you to R as a programming language and statistical package (using the RStudio user interface)

The class

Aims:

1. Introduce you to R as a programming language and statistical package (using the RStudio user interface)
2. Implement basic data analysis tasks in R, including data management, visualization, and regression analysis

The class

Aims:

1. Introduce you to R as a programming language and statistical package (using the RStudio user interface)
2. Implement basic data analysis tasks in R, including data management, visualization, and regression analysis
3. Reproduce existing Political Science research

Schedule

1. Getting set up, 26 March, 18-20h (Aula 10, Sede Centrale)
2. Data management, 2 April, 18-20h (Aula 2, Palazzo Boyl)
3. Visualization, 9 April, 18-20h (Aula 10, Sede Centrale)
4. R programming, 23 April, 18-20h (Aula 5, Sede Centrale)
5. Exploring data, 7 May, 18-20h (Aula 3, Palazzo Boyl)
6. Linear Regression, 14 May, 9-11h (Aula 3, Sede Centrale)
7. Replication I, 14 May, 11-13h (Aula 3, Sede Centrale)
8. Logistic regression, 21 May, 9-11h (Aula 3, Palazzo Boyl)
9. Replication II, 21 May, 11-13h (Aula 3, Palazzo Boyl)
10. Conclusion, 28 May, 18-20h (Aula 10, Sede Centrale)



Today's class

1. RStudio



Today's class

1. RStudio
2. Basic notions in R: Objects, functions, packages



Today's class

1. RStudio
2. Basic notions in R: Objects, functions, packages
3. Some statistical notions



Today's class

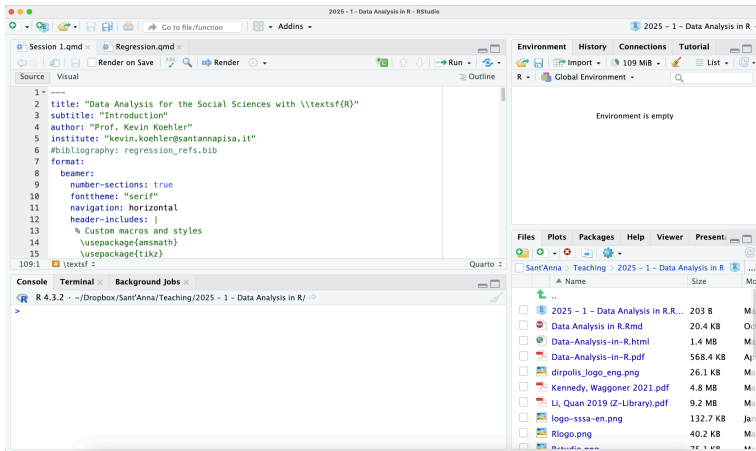
1. RStudio
2. Basic notions in R: Objects, functions, packages
3. Some statistical notions
4. Exercises



RStudio



RStudio



First, create a **project**. Projects are useful for keeping all files in the same place.

1. Create a directory on your computer where you want to save all files related to this class
2. Go to **> File > New Project** in the R menu and then select “Existing directory”
3. Navigate to the folder and name and create the project



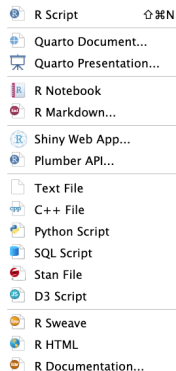
The folder you created is now also your **working directory**. You can see the directory at the top of the console. You can also type:

```
getwd()
```



Basic notions in R

Types of R files



R can do many different things, not just statistical analysis.

Consequently, there are many different file types in RStudio:

- ▶ R scripts for coding
- ▶ Quarto documents and presentations
- ▶ R Notebook and R Markdown
- ▶ Shiny Apps
- ▶ Plumber API
- ▶ Files in other languages (C++, Python, SQL...)



R scripts

- ▶ Go ahead and open a new R script.
- ▶ Type `print("Hello world")`
- ▶ With your cursor in the line with the command, press **Ctrl + Enter** (Windows) or **Command + Enter** (Mac)
- ▶ The code is executed and the results printed in the Console



objects



Vectors

Vectors are the most basic data structure. They hold a series of numeric or character values and are created with the `c()` function (the `c` in the function stands for concatenate):

```
c(1,2,3,4,5)
```

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

```
c("a","b","c","d","e")
```

```
[1] "a" "b" "c" "d" "e"
```



Matrices and data frames

```
matrix(1:6, nrow = 2, ncol = 3)
```

	[,1]	[,2]	[,3]
[1,]	1	3	5
[2,]	2	4	6

```
data.frame(name = c("A","B"),  
            age = c(24,56))
```

	name	age
1	A	24
2	B	56



The environment

The **Environment** is where R stores objects for the duration of a session. You can assign an object to the environment by typing:

```
data <- data.frame(name = c("A","B"),  
                   age = c(24,56))
```

(you could replace the <- with =, but I recommend getting used to <-)

After running this code, you should have an **object** called “data” in your environment. This object contains two columns (name and age) with two rows each. We call the columns **variables** and the rows **observations**.

You can click on the object to see what it contains.

Working with data frames

```
data
```

```
  name age  
1    A  24  
2    B  56
```

```
data$name
```

```
[1] "A" "B"
```

```
data$age[1]
```

```
[1] 24
```

```
data$age[data$name=="A"]
```

```
[1] 24
```

Functions

R works with **functions**. A function takes an **object** (or multiple **objects**) as input and does something with it.

Examples include:

- ▶ `print("Hello world")` prints “Hello world” to the console
- ▶ `c(1,2,3,4)` creates a numerical vector with 1,2,3,4 as elements
- ▶ `getwd()` returns the active working directory
- ▶ `help(print)` returns the help file for the `print()` function
- ▶ `lm(x~y)` performs a linear regression of x on y

Functions in R are words followed by brackets. You always need to close the brackets, otherwise your code will not run.

User-defined functions

You can write your own functions in R. Here is a function which takes a number as an argument and tells you whether the number is greater than 5:

```
greater5 <- function(x) {  
  if (!is.numeric(x)) {  
    stop(paste0("Argument must be numeric.\n",  
                "You provided an object of class: ",  
                class(x)[1],  
                ". You moron."))  
  } # check if input is numeric, return error if not  
  result <- ifelse(x > 5,  
                   paste(x, "is greater than 5"),  
                   paste(x, "is not greater than 5"))  
  return(result) # Return results  
}
```



Packages and CRAN

Functions are part of packages. Your version of R comes with base R, but there are many other packages.

We will use the **tidyverse** family of packages. You can install packages by typing `install.packages("tidyverse")` in the Console. This will download the package and save it on your machine. You need to do this only once.

To use specific packages, you need to load them in the beginning of your R session. It is good practice to include all packages needed to run your code in the beginning of your R script. Packages are loaded typing `library(tidyverse)`.

Some statistical notions



Recap on variable types

1. **Nominal Scale:**

- ▶ Categories without a specific order (e.g., gender, color).

2. **Ordinal Scale:**

- ▶ Categories with a defined order but unequal intervals (e.g., rankings, satisfaction ratings).

3. **Interval Scale:**

- ▶ Numeric scales with equal intervals but no true zero (e.g., temperature in Celsius).

4. **Ratio Scale:**

- ▶ Numeric scales with equal intervals and a true zero (e.g., height, weight, age).



How can we describe the typical value for
each of these scales?



Measures of Central Tendency

- ▶ **Mean:** The average of a data set, calculated by summing all values and dividing by the number of values.



Measures of Central Tendency

- ▶ **Mean:** The average of a data set, calculated by summing all values and dividing by the number of values.
 - ▶ Mean = $\frac{\sum_{i=1}^n x_i}{n}$, where x_i represents each data point, and n is the number of data points.



Measures of Central Tendency

- ▶ **Mean:** The average of a data set, calculated by summing all values and dividing by the number of values.
 - ▶ Mean = $\frac{\sum_{i=1}^n x_i}{n}$, where x_i represents each data point, and n is the number of data points.
- ▶ **Median:** The middle value when data is ordered from lowest to highest; 50% of values fall below it.



Measures of Central Tendency

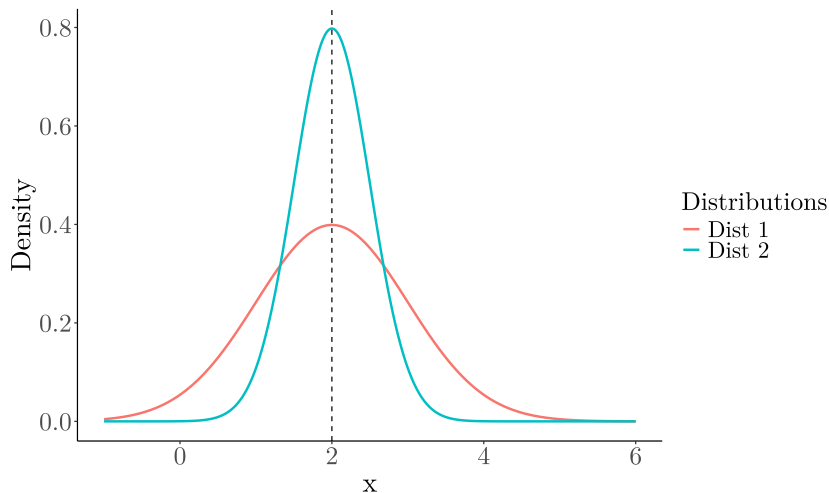
- ▶ **Mean:** The average of a data set, calculated by summing all values and dividing by the number of values.
 - ▶ Mean = $\frac{\sum_{i=1}^n x_i}{n}$, where x_i represents each data point, and n is the number of data points.
- ▶ **Median:** The middle value when data is ordered from lowest to highest; 50% of values fall below it.
- ▶ **Mode:** The value that occurs most frequently in a data set.



How can we describe variation around a central value?



How can we describe variation around a central value?



Measures of Dispersion

- ▶ **Range:** $\text{Range} = \text{Max}(x) - \text{Min}(x)$. The difference between the maximum and minimum values in the data set.



Measures of Dispersion

- ▶ **Range:** $\text{Range} = \text{Max}(x) - \text{Min}(x)$. The difference between the maximum and minimum values in the data set.
- ▶ **Variance:** $\text{Variance} = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$, where x_i represents each data point, μ is the mean, and n is the number of data points.



Measures of Dispersion

- ▶ **Range:** $\text{Range} = \text{Max}(x) - \text{Min}(x)$. The difference between the maximum and minimum values in the data set.
- ▶ **Variance:** $\text{Variance} = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$, where x_i represents each data point, μ is the mean, and n is the number of data points.
- ▶ **Standard Deviation:** $\text{SD} = \sqrt{\text{Variance}}$. The square root of the variance, which gives the spread of data in the same units as the original data.



Measures of Dispersion

- ▶ **Range:** $\text{Range} = \text{Max}(x) - \text{Min}(x)$. The difference between the maximum and minimum values in the data set.
- ▶ **Variance:** $\text{Variance} = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$, where x_i represents each data point, μ is the mean, and n is the number of data points.
- ▶ **Standard Deviation:** $\text{SD} = \sqrt{\text{Variance}}$. The square root of the variance, which gives the spread of data in the same units as the original data.
- ▶ **Interquartile Range (IQR):** $\text{IQR} = Q3 - Q1$, where $Q3$ is the third quartile and $Q1$ is the first quartile, representing the middle 50% of the data.



Exercises



GitHub repository

There is a GitHub repository for this class where I will share materials with you. You can reach it here:

<https://github.com/KevinKoehlerSSSA/Intro-to-R>

Exercises

1. Read the data into R, save it in your environment as `tun22`.
2. Consult the codebook to understand what you are looking at
3. Write code to calculate:
 - 3.1 The typical age of all respondents
 - 3.2 The typical age of respondents who have voted for Kais Saied in the first round of the presidential elections
4. Describe how much Saied voters differ from each other in terms of the levels of education. Which measure would you use? Why?
 - 4.1 Write the appropriate code



Additional exercises

1. What is the vote percentage of Kais Saied in the first and second round?
2. Are first round Saied voters significantly younger or older than respondents overall?
3. Are male respondents significantly more likely to have voted Saied in the first round than female respondents?

