

Introduction to R

Quantitative Methods I

Malo Jan

2024-09-12

Table of contents

Course overview	4
Course structure	4
Course validation	4
Requirements	5
Course material	5
Help	5
Ressources	6
Setting up R and Rstudio	7
Install R	7
Install Rstudio	7
Set preferences and check if everything works	8
Installation issues	9
Mac	9
Windows	9
I Getting started with R and Rstudio	10
1 Rstudio, scripts, quarto	11
1.1 How Rstudio is organized	11
1.2 Scripts	12
1.3 Quarto	12
2 R basics	14
2.1 R as calculator	14
2.2 Objects	15
2.3 Vectors	16
2.3.1 Characters vectors	17
2.3.2 Logical vectors	18
2.4 Functions	19
2.5 Missing values	20
2.6 Packages and libraries	21
2.7 Dataframes and tibbles	22
2.8 Import and write data	24
2.8.1 Getting data	24

2.8.2	Read data into R	24
2.8.3	Functions to import data in R	26
2.8.4	Export data	27
2.9	Going further	27
2.9.1	Naming things	27
2.9.2	Packages	27
2.9.3	Indexing	28
2.9.4	More on vectors and functions	29
3	Getting help and dealing with errors	33
3.1	R documentation	33
3.2	Cheatsheets	33
3.3	Online ressources	33
3.4	AI is your friend	34
3.5	Most common errors	34
II	Manipulating and describing data	38
III	Visualizing data	39
IV	Testing relationships	40
V	Regression analysis	41
VI	Multivariate analysis	42

Course overview

This website contains all the material for the lab sessions of the Quantitative Methods class for the Fall 2025 semester in the [research master's program in Political Science](#) at Sciences Po Paris. The class complements Jan Rovny's lecture on [Quantitative Methods I](#).

The course provides you with the fundamentals, resources, and motivation to further your learning independently, prepare for the next semester [Luis Sattelmayer's R sessions](#), and apply quantitative methods in your future research.

Course structure

Session	Description	Date
Session 1	Getting started with R and Rstudio	05/09, 12/09
Session 2	Manipulating and describing data	19/09, 26/09
Session 3	Visualizing data	03/10, 10/10
Session 4	Testing relationships	17/10, 24/10
Session 5	Correlation and simple linear regression	07/11, 14/11
Session 6	Multivariate analysis	21/11, 28/11

Course validation

Learning programming is fundamentally about practice. It involves trying, encountering challenges, and solving them. The course assessment is structured around a series of exercises that will constitute 30% of your final grade in the Quantitative Methods class. You will be required to apply the code and concepts covered in class to new problems and datasets.

Throughout the semester, you will complete four individual exercises between sessions. For each exercise, you will perform a series of operations in R and submit your work. Each exercise will be graded on a scale from 0 to 5 points.

Second, at the end of the semester, you will complete a bigger group exercise with one other person. This exercise will involve integrating content from throughout the semester and will be completed within a two-week timeframe. It will be graded on a scale from 0 to 20 points.

Assignments	Description	Due date	Weight
Exercise 1	Individual	7 days after Session 1	5%
Exercise 2	Individual	7 days after Session 2	5%
Exercise 3	Individual	7 days after Session 3	5%
Exercise 4	Individual	7 days after Session 4	5%
Exercise 5	Group	14 days after session 6	10%

Requirements

This class does not require any prior programming or statistical experience and is designed for complete beginners. However, it does require some basic knowledge of how to use a computer. You should be able to navigate your computer's file system, create and move folders, and download and save files. Specifically, each class's content will be provided as a zip file, which you will need to download, unzip, and move to a folder on your computer. If you are not familiar with these operations, please refer to this [tutorial](#).

This class requires you to bring a laptop to each session with R and RStudio installed by the start of the course. The 'Setting up R and RStudio' page on this website provides a step-by-step guide for installation. If you encounter any issues during the setup or are unable to bring a laptop to class, please let me know.

Course material

Before each session, you will be provided with the class material, which you should download to your laptop to follow along. Additionally, all class activities and resources will be available on this website for future reference. You will also need to submit your exercises on the course's Moodle page.

Help

If you have any questions regarding the course, need help, or are looking for additional resources, please do not hesitate to contact me via [email](#). I will be happy to assist you and will try to reply as quickly as possible.

Ressources

This course does not have a required textbook or mandatory readings. However, if you wish to deepen your understanding of the content covered, I recommend familiarizing yourself with the following resources:

- [R for Data Science](#), this is THE R classic by Hadley Wickham, you should definitely take a look to better understand what we cover in this course.
- [Telling stories with data](#) by Rohan Alexander : one of my favorite book on data science with R. (a bit more advanced)
- [Computational analysis of communication](#) by van Atteveldt et al.
- [Computational Thinking for Social Scientists](#) by Jae Yeon Kim
- [Introduction to data science](#) by Rafael Irizarry

There are also many other introductions to R available online, each offering different approaches to teaching the same concepts. I recommend the following:

- [Introduction to R](#), by Felix Lennert
- [Introduction to R](#) by Alex Douglas et al.

If you want ressources in french, these are the two most comprehensive introduction you will find :

- [Introduction à R et au tidyverse](#) by Julien Barnier
- [Guide pour l'analyse de données d'enquêtes avec R](#) by Joseph Larmarange

Setting up R and Rstudio

The goal of this course is to apply what you learn in the introduction to quantitative methods' course by learning **programming**. This involves giving instructions to your computer to perform operations on quantitative data using code. The course focuses on coding with the programming language R, using the software RStudio.

To get the first session off to a good start and avoid wasting too much time setting up, **you must bring a laptop with R and Rstudio installed**. Instructions are provided below. This should be fairly quick. If you have any issues to install these programs correctly, **please contact me before the first class**.

R and RStudio are two different things. R is a programming language and software widely used for statistical analysis. While it is possible to use the default software directly, many (if not most) prefer to use RStudio instead. RStudio is an integrated development environment (IDE) that provides a much more functional and user-friendly interface to interact with R, which is why we choose to use it. However, we still need to install R before installing RStudio, as the latter relies on the former to function properly.

Install R

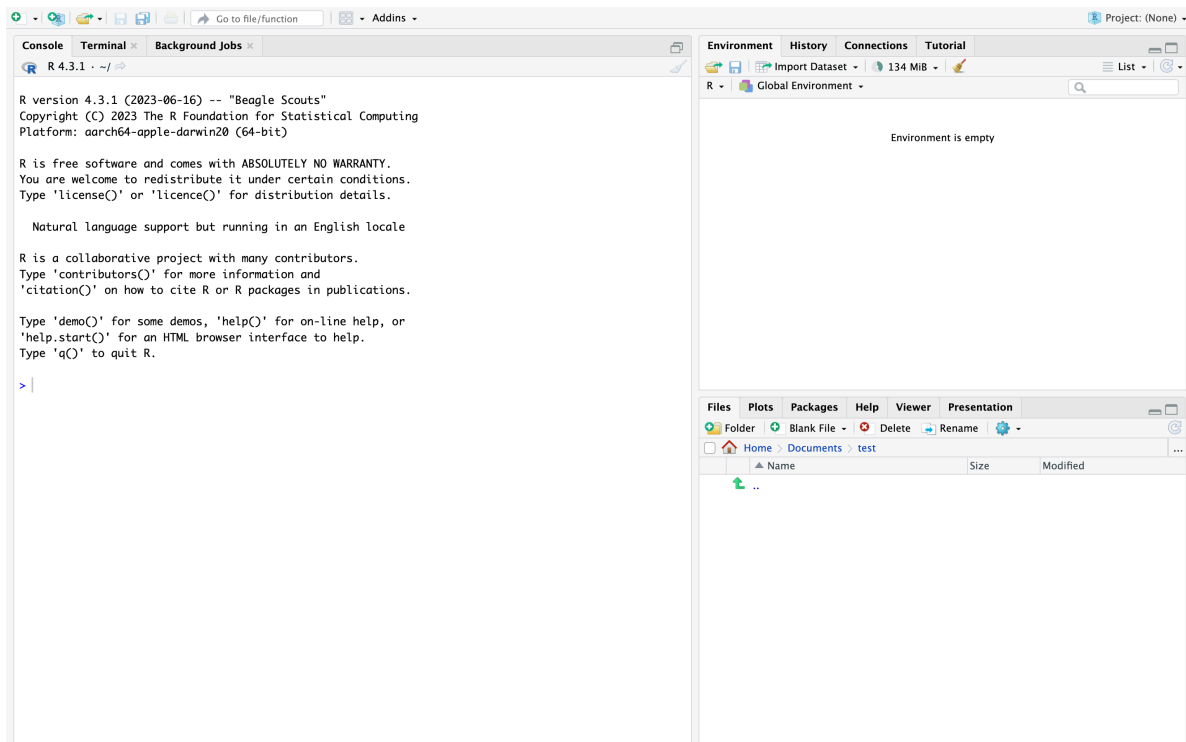
You can install R from [CRAN](#) (The comprehensive R Archive Network). On CRAN, you will find links to download the version of R you need for your operating system (Windows, Mac or Linux). Once the download is complete, you need to execute the installer.

Install Rstudio

Once R is set up, you can install Rstudio. For this, go on this [webpage](#), download Rstudio and follow the instructions.

Set preferences and check if everything works

Once the installation of both R and Rstudio is complete, you should open Rstudio. If everything worked, something very similar to the screenshot below should appear. If it doesn't open, restart your laptop and try again.



RStudio has some default settings that are worth changing to adopt best practices. Please do the following:

- Tools /Global options/RGeneral : Save workspace to never and uncheck the box **Restore .Rdata into workspace at startup**.

While R provides a series of basic commands for data manipulation, many of the functionalities we will use come from **packages** that need to be installed and loaded via RStudio. The most well-known package we will use is called **tidyverse**. To check if your installation is successful and you are ready to start the course, copy-paste the following code into the console pane and press enter. It might take a few minutes, and at the end, you should see the same message as the one below.

```
if (!require(tidyverse)) {install.packages("tidyverse")}
```

Loading required package: tidyverse


```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0.9000 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2    3.5.1      v tibble     3.2.1
v lubridate  1.9.3      v tidyr      1.3.1
v purrr      1.0.2
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

Installation issues

If you do not see the message above, welcome in the world of troubleshooting! Depending on your operating system (MacOS, Windows, Linux) and the configuration of your computer, the installation process may not go as smoothly as expected. Here are some common issues and how to solve them:

Mac

If a message appears saying that the package cannot be found, try the following :

```
if (!require(tidyverse)) {devtools::install_github("hadley/tidyverse")}
```

Windows

If a message appears saying that the package cannot be found, try to install Rtools by copying and pasting the following code into the console pane.

```
if (Sys.info()["sysname"] == "Windows") {
  install.packages("Rtools")
}
```

If a message appears that Rtools is not available, download it from [here](#). Then restart your computer and try to install the tidyverse package again.

If all of this does not work or you have different types of issues, please contact me.

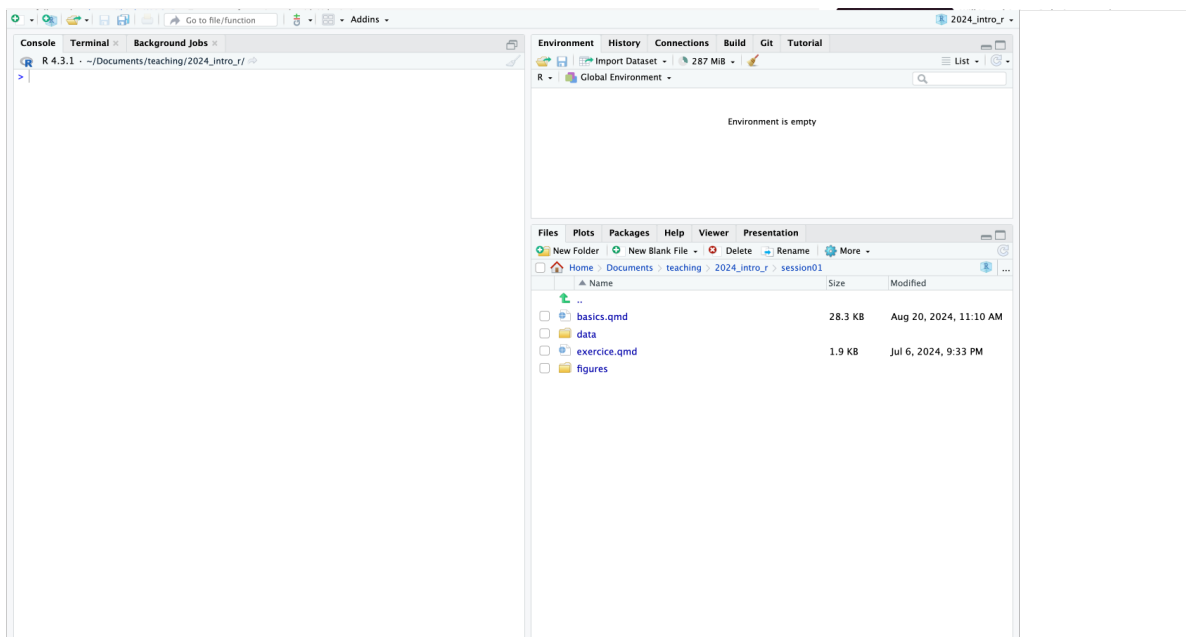
Part I

Getting started with R and Rstudio

1 Rstudio, scripts, quarto

1.1 How Rstudio is organized

To interact with the language R, we use the interface of Rstudio. When you open it on your laptop, you will probably see something similar to the screenshot below.



On the upper right panel, you have what we call the *environment*. At this stage, you should see a message indicating that your environment is empty. However, when we use data in RStudio, what we manipulate will appear there. For instance, if I import a dataset, I will see in that pane that my dataset has been imported.

On the bottom left panel, you see what we call our working directory. This is where the files, such as the data, that we want to use are stored on our laptop. It is essential to be aware that if you remove a file here from RStudio, it will also be deleted from your laptop. Exercise caution when managing files in this directory to avoid unintentional deletions.

On the left side, you have the console where you can type commands, indicated by the > sign. If you type something here and then click on **Enter**, you will see the result. So, if you type

2+2, it will produce the output 4 just below. Whatever code you produce, the output will appear in the console.

1.2 Scripts

However, when we interact with R, we rarely use the console to type code. Instead, we mostly use *scripts*, which are specific types of documents where you can write code, save it, and reuse it later on. To open a new script, you can go into the menu bar, click on **File > New File > Rscript** (you can also use the shortcut **Ctrl+Shift+N** on Windows/Linux or **Cmd+Shift+N** on macOS). This will open a new document, known as an R script, in the editor on the upper left side of the interface. Here, you can start writing and saving your R code for analysis and projects. You can type again 2+2. To run a line on a Rscript, place the cursor on the line of code you want to run and press **Ctrl + Enter** (on Windows/Linux) or **Cmd + Enter** (on macOS). Just above the script in RStudio, you will also find a **Run** button that does the same thing. Once you run that code, you should see the output printed in the console. Note that in R scripts, you can write text and comment your code with the **#**. To save your script, simply click on the small disk icon in the editor's toolbar and choose a name for the file (**Ctrl + S** on Windows/Linux or **Cmd + S** on macOS). You will now see it in the files pane at the bottom left of your screen.

```
2+2 # This is my first code
```

```
[1] 4
```

1.3 Quarto

While many users opt for R scripts for coding, I won't extensively employ them in this class. Instead, I'll be using Quarto documents. You will also be using Quarto as I'll be requesting you to submit your assignments using this format. Quarto constitutes a method of cohesively editing code and text within a single document, a practice termed *literate programming*. It enables you to compose code, generate output, conduct analyses in the same document, format the content, and subsequently export it to diverse formats like Word documents, PDFs and html. Additionally, Quarto facilitates the creation of slides, dashboards, books, and websites. Personally, I employ Quarto to craft the course materials.

To create a quarto document, click on **File > New File > Quarto Document** in the menu bar. You will see different options appear but for now, uncheck **use virtual markdown editor** and click on **create an empty document**. There are three main differences with an R script :

- **YAML** : when you open a quarto document you will have to fill the top matter, called a call a **YAML** which is separated by `---` where you can write different informations and options such as the title, the author, subtitle, abstract, the data and many other things. You can also choose the format of the output by writing **format:** and choosing among **html**, **docx** or **pdf**.
- You can type text without the `#` as if you where typing in any text editor
- To write code, you need to create what we call a chunk. For this, either go on the menu bar : **Code > Insert Chunk** or **option + Cmd + i** in macOS, **Ctrl+Alt+i** on Windows/Linux. You can then directly type code in that chunk.

To generate a pdf/html/word document, you need to click on **Render** (**Cmd + Shift + K** in macOS, {< kbd win=Shift-Ctrl-P >}). A new file will appear in your files pane with the output.

2 R basics

2.1 R as calculator

First of all, R is a fancy calculator that can be used to perform fundamental arithmetic operations.

```
3+7+10 # Addition
```

```
[1] 20
```

```
4-5 # Substraction
```

```
[1] -1
```

```
3*9*10 # Multiplication
```

```
[1] 270
```

```
2/6 # Division
```

```
[1] 0.3333333
```

```
2^2 # Exponentiation
```

```
[1] 4
```

```
(2+2)-(4*4)/2^2 # Mix of operations
```

```
[1] 0
```

2.2 Objects

When you run code in R, the results are shown in the console. However, you cannot directly reuse these results in further operations, which is what we want to do. To address this, we use **objects**. Objects in R act as containers that store values, allowing you to keep information for later use. To create an object in R, you use the **assignment operator** `<-` .

```
my_object <- 2
my_object
```

```
[1] 2
```

Let's consider an example from the most recent French general election. The results led to a highly fragmented parliament, with 11 different parliamentary groups and no single party or coalition able to form an absolute majority on its own. This situation raises the question of which coalitions could be formed to achieve a majority and pass legislation. The code below uses objects to store the number of seats obtained by each parliamentary group.

```
# Left-wing groups
communists_seats <- 17
lfi_seats <- 72
greens_seats <- 38
socialists_seats <- 66

# Macron's party and allies
renaissance_seats <- 99 # This is Macron's party
modem_seats <- 36       # MODEM, a centrist party
horizon_seats <- 31     # Horizon, party of the former PM Édouard Philippe

# Right-wing groups
conservatives_seats <- 47 # "Les Républicains"

# Far-right groups
rn_seats <- 126          # National Rally (RN)
ciotti_seats <- 16       # Former conservatives who allied with Le Pen's party

# Others
liot_seats <- 22         # A mix of some independent centrists and regionalists
none_seats <- 7          # Non-affiliated MPs
```

After executing these lines, you should be able to see the objects created and their values in the environment pane located in the upper-right section of RStudio. Once these objects

created, it is possible to perform operations on them. For instance, it is possible to compute the total number of seats of actual existing coalitions by summing the number of seats of the different parties that compose them and saving them in new objects. Below, I create a new object `left_seats` that stores the total number of seats obtained by the left-wing parties by summing the different objects that store the number of seats obtained by the left-wing parties.

```
left_seats <- communists_seats + lfi_seats + greens_seats + socialists_seats
left_seats
```

```
[1] 193
```

i Your turn !

Create two new objects `macron_seats` and `far_right_seats` that store the total number of seats obtained by Macron's allies and far-right parties respectively.

Solution. Click to expand!

Solution:

```
macron_seats <- renaissance_seats + modem_seats + horizon_seats
far_right_seats <- rn_seats + ciotti_seats
```

2.3 Vectors

The objects we used so far contained only **one** numeric value. However, what we mostly manipulate in R are **vectors**, which are sequences of different values on which we can perform operations. Vectors can be of different types (eg : numeric, character, logical, date) but they have to be of the same type. For instance, a numeric vector is a sequence of different numbers and a character vector is a sequence of different strings. Vectors are also unidimensionals which mean they contains only one sequence of values and not several such as matrices do.

We can generate vectors with `c()` which stands for “concatenate”. For instance, here, I create a vector containing the values of the seats obtained by different coalitions in the election. As a result, the vector `coalition_seats` contains the number of seats obtained by the left, Macron's party and allies, the far-right and the conservatives.

```
coalition_seats <- c(left_seats, macron_seats, far_right_seats, conservatives_seats)
coalition_seats
```

```
[1] 193 166 142 47
```


We use vectors to store different values because it is possible to perform the same operation on all the values of a vector at once. Let's say we want to know the number of seats that an existing coalition would need to have an absolute majority in the parliament. We can do this by creating a new object, `majority`, that stores the number of seats needed for an absolute majority, and then subtracting the number of seats obtained by the different coalitions from this value. The result will be a vector showing the number of seats each coalition needs to reach an absolute majority, which we could also store in a new object if desired.

```
majority <- 577/2 + 0.5 # Half of the number of seats (577) + 0.5 to round up (there is no h  
majority
```

```
[1] 289
```

```
majority - coalition_seats
```

```
[1] 96 123 147 242
```

2.3.1 Characters vectors

So far, we have only used numerical vectors, which consist of numbers. However, we can also create character vectors, which are made up of strings enclosed in quotes (either single ' or double "). For example, we can create a vector containing the names of different parliamentary leaders.

```
left_leaders <- c("Chassaigne", "Chatelain", "Panot", "Vallaud") # Create a vector of left-w  
far_right_leaders <- c("Ciotti", "Le Pen") # Create a vector of far-right leaders
```

As for other vectors, you can combine them in a single vector which will return a vector with all the leaders' names.

```
leaders <- c(left_leaders, far_right_leaders)  
leaders
```

```
[1] "Chassaigne" "Chatelain" "Panot"      "Vallaud"    "Ciotti"  
[6] "Le Pen"
```

2.3.2 Logical vectors

Another type of vector in R is the logical vector, which consists of Boolean values: `TRUE` or `FALSE`. Logical vectors are useful for evaluating conditions. They can be used to check for errors in data or to filter variables based on specific criteria.

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

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

For instance, we could check whether there are no mistakes and confirm that Le Pen is not a leader of the left coalition. This means verifying that no value in the `left_leaders` vector is equal to “Le Pen”.

```
left_leaders == "Le Pen"
```

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

```
far_right_leaders == "Le Pen"
```

```
[1] FALSE TRUE
```

The evaluation of conditions can also be used to compare numeric values. For instance, we can check whether the left coalition has more seats than the far-right coalition.

```
left_seats > far_right_seats
```

```
[1] TRUE
```

Your turn !

Create a vector that contains the names of the leaders of the Macron’s party : Attal, Fesneau, Marcangeli. Then, add this vector to the `leaders` vector and store the result in a new object `all_leaders`.

Solution. Click to expand!

Solution:

```

macron_leaders <- c("Attal", "Fesneau", "Marcangeli")
all_leaders <- c(leaders, macron_leaders)

all_leaders[c(5,6,7)]

[1] "Ciotti" "Le Pen" "Attal"

all_leaders[all_leaders %in% c("Attal", "Fesneau", "Marcangeli")]

[1] "Attal"      "Fesneau"    "Marcangeli"

```

2.4 Functions

To manipulate vectors and conduct operations on them, we use **functions**. A function is a reusable block of code that performs a specific task, it takes several input values called *arguments* and produce an *output*.

Let's say you want to know how many seats parliamentary groups have on average in the French parliament. You could calculate the sum of the seats and dividing them by their number.

```

(communists_seats + greens_seats + socialists_seats + lfi_seats + renaissance_seats + modern_seats) / length(communists_seats)

[1] 48.08333

```

But you could also just the `mean()` function that exists in R and that takes a vector of numbers as argument.

```

parl_seats <- c(communists_seats, greens_seats, socialists_seats, lfi_seats, renaissance_seats, modern_seats)

mean(parl_seats) # Compute the mean of the vector parl_seats

```

```

[1] 48.08333

```

It is also easy to use functions to compute the maximum, minimum, and sum of a vector of numbers.

```
max(parl_seats) # Compute the maximum of the vector parl_seats
```

```
[1] 126
```

```
min(parl_seats) # Compute the minimum of the vector parl_seats
```

```
[1] 7
```

```
sum(parl_seats) # Compute the sum of the vector parl_seats
```

```
[1] 577
```

2.5 Missing values

In R, a missing value is represented by the symbol `NA`, which stands for “Not Available.” Missing values can arise for a variety of reasons, such as data not being observed or recorded, errors in data collection, or intentional omissions. Understanding and handling missing values is crucial because they can influence the results of your analysis or even cause some functions to return errors. For instance, imagine I haven’t found any information about the number of seats one parliamentary group got, but I want to retain this information in my vector. So, I add an `NA` to it.

```
parl_seats <- c(parl_seats, NA)
parl_seats
```

```
[1] 17 38 66 72 99 36 31 47 16 126 7 22 NA
```

When analyzing data, it’s not uncommon to encounter `NA` values, and it’s important to be aware of them. To check if a vector contains `NA` values, you can use the `is.na()` function. This function returns a logical vector indicating whether each value is `NA` (`TRUE`) or not (`FALSE`).

```
is.na(parl_seats) # Check which values of a vector are NAs
```

```
[1] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
[13] TRUE
```

This is important because certain functions will not operate properly if there are NA values in your data. For instance, the `mean()` function will return NA if the data contains any NA values

```
mean(parl_seats)
```

```
[1] NA
```

To deal with NA, the `mean()` function has a `na.rm`

```
mean(parl_seats, na.rm = TRUE) # Remove NA before computing the mean
```

```
[1] 48.08333
```

2.6 Packages and libraries

The functions we've discussed so far, such as `sum()` and `mean()`, come from base R. These are pre-loaded functions available immediately upon starting R. However, many functions you'll encounter aren't part of base R but instead belong to specific packages that individuals or groups have developed. You can think of packages as collections of functions crafted to simplify certain tasks or to introduce new capabilities to R. For example, there's the `tidyverse` package, which I asked you to install before the class

To install a package in R, you can use the `install.packages()` function, passing the name of the package in quotation marks (either single or double). I recommend doing this installation in the console since you don't need to save this step; it's a one-time action. However, every time you start your script or Quarto document, you'll need to load the package. To do this, use the `library()` function, providing the package name as an argument, but without the quotation marks.

```
library(tidyverse)
```

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0.9000 --
```

```
v dplyr      1.1.4      v readr      2.1.5
```

```
v forcats   1.0.0      v stringr    1.5.1
```

```
v ggplot2    3.5.1      v tibble     3.2.1
```

```
v lubridate  1.9.3      v tidyr      1.3.1
```

```
v purrr      1.0.2
```

```
-- Conflicts ----- tidyverse_conflicts() --
```

```
x dplyr::filter() masks stats::filter()
```

```
x dplyr::lag()     masks stats::lag()
```

```
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

2.7 Dataframes and tibbles

When processing data, we primarily interact with vectors through the manipulation of **dataframes** in R. Dataframes are two-dimensional structures that contain rows and columns. Each column is a vector, and each row is an observation. Dataframes are the primary data structure used for data manipulation, computation, and visualization in R. In this class, we'll work with a specific type of dataframe that comes from the **tidyverse** package called a **tibble**. Tibbles are a modern and enhanced version of dataframes that make them easier to print and manipulate.

However, we primarily interact with vectors through the manipulation of dataframes in R. Dataframes are composed of combinations of vectors, which can vary in types. Dataframes are two-dimensional, with columns (or variables) and rows (or observations). This is what we use for manipulating data, computing statistics and visualization. In this class, we will work with a specific form of dataframe coming from the **tidyverse** packages that is called a **tibble**. Tibbles make dataframes easier to print and manipulate.

To understand what dataframes look like, let us continue with the results of the spanish elections. I manually create a tibble with the `tibble()` function with different variables about different parties, their seats, their vote share and their candidate.

```
left_coalition <- tibble(  
  party = c("Communists", "Greens", "LFI", "Socialists"),  
  leader = c("Chassaigne", "Chatelain", "Panot", "Vallaud"),  
  seats_2024 = c(17, 38, 72, 66),  
  seats_2022 = c(22, 21, 75, 31)  
)
```

You see now that we have a new object in our Environment Pane with 3 observations and 5 variables. If we want to access only one variable (one vector) of that dataframe, we use the `$` sign. This will return a vector of the values of this variable. You can also get the same result by indicating the position of the column inside `[[]]`.

```
left_coalition$party # Select the party variable
```

```
[1] "Communists" "Greens"      "LFI"         "Socialists"
```

Different functions are available to get an idea of the informations and shape of the dataframe, which are useful when we load an unknown dataset and we want to understand its structure, what are the observations and variables.

```
head(left_coalition, 1) # Return x first rows of an object
```

```
# A tibble: 1 x 4
  party      leader seats_2024 seats_2022
  <chr>      <chr>      <dbl>      <dbl>
1 Communists Chassaigne      17        22
```

```
tail(left_coalition, 2) # Return x last rows of an object
```

```
# A tibble: 2 x 4
  party      leader seats_2024 seats_2022
  <chr>      <chr>      <dbl>      <dbl>
1 LFI        Panot        72        75
2 Socialists Vallaud      66        31
```

```
dplyr::glimpse(left_coalition) # Get a glimpse of your data
```

```
Rows: 4
Columns: 4
$ party      <chr> "Communists", "Greens", "LFI", "Socialists"
$ leader     <chr> "Chassaigne", "Chatelain", "Panot", "Vallaud"
$ seats_2024 <dbl> 17, 38, 72, 66
$ seats_2022 <dbl> 22, 21, 75, 31
```

```
colnames(left_coalition) # Retrieve column names of the dataframe
```

```
[1] "party"      "leader"     "seats_2024" "seats_2022"
```

```
summary(left_coalition) # Return a summary of the variables
```

party	leader	seats_2024	seats_2022
Length:4	Length:4	Min. :17.00	Min. :21.00
Class :character	Class :character	1st Qu.:32.75	1st Qu.:21.75
Mode :character	Mode :character	Median :52.00	Median :26.50
		Mean :48.25	Mean :37.25
		3rd Qu.:67.50	3rd Qu.:42.00
		Max. :72.00	Max. :75.00

You can also create new variables based on the existing ones. Here I create a new variable called `seats_share` by calculating the share of seats each party has in the parliament (dividing their seats by the total number of seats and multiplying by 100).

```
left_coalition$seats_share <- left_coalition$seats_2024/577
```

i Your turn !

Create a new variable called `seats_change` that calculates the difference of seats between 2024 and 2022 for each party.

Solution. Click to expand!

Solution:

```
left_coalition$seats_change <- left_coalition$seats_2024 - left_coalition$seats_2022
```

2.8 Import and write data

2.8.1 Getting data

Up until now, we have been creating data manually for demonstration purposes. As we move forward, we will focus on analyzing real data. Every data analysis project starts with acquiring data. You can generate your own data through surveys, web scraping, or manual data coding (*primary data sources*). However, there are also many pre-existing datasets available for use (*secondary data sources*). These datasets are often provided by researchers, governments, NGOs, companies, international organizations, and more. I highly recommend checking out this [list](#) of political datasets curated by Erik Gahner Larsen. Throughout this course, we will use some of these established datasets in political science.

2.8.2 Read data into R

2.8.2.1 File formats, paths and R projects

To analyze data in R, we first need to import it. While this might sound simple, it can be challenging for beginners. To read a file in R, we need to know two important things: the **file format** and the **path**.

Data comes in various file formats, which are standardized ways of storing and organizing data in digital files. These formats dictate how information is encoded and structured, allowing different software programs to understand and interpret the data correctly. The most common format for data is `.csv` (comma-separated values). In political science, you will also encounter

Stata (.dta) and SPSS (.sav) files. R uses different functions to read files depending on their format.

Secondly, R needs to know where the data is located on your system. A “path” shows the position of a file or folder within your file system, detailing the series of directories and subdirectories leading to the file. There are two types of paths:

- The **absolute path** gives the complete location of a file or directory, beginning at the root of the file system. Examples include: /home/user/documents/myfile.txt for Unix-like systems, and C:\Users\user\Documents\myfile.txt for Windows.
- A **relative path** indicates the location of a file or directory in relation to the current working directory. For instance, data/mydata.csv points to a file named mydata.csv in the data subdirectory of the present directory.

The **working directory** refers to the directory where R is currently operating. If you access files in R without providing an absolute path, it defaults to searching within this working directory.

You can see your current working directory in R using the `getwd()` function. To set a new working directory, use the `setwd()` function, specifying the desired path as its argument.

```
getwd()
```

```
[1] "/Users/malo/Documents/teaching/2024_intro_r/session01"
```

Please note that using absolute paths in your code is considered a bad practice because it can make your code less usable for others. Instead, I recommend using R projects. An R project is a dedicated workspace for your R work, where you keep all your files, data, scripts, and output documents together. When you open an R project, it sets everything up so that your files are easy to find and your work is easier to share and reproduce.

From the top left corner of RStudio, click on **File** and then select **New Project**. You’ll then be given three options: to create a project in an existing directory, to create one in a new directory, or to check out a project from a version control repository like Git.

If you choose an existing directory, navigate to that directory. If you opt for a new directory, you’ll need to name your project and decide its save location on your computer. Once you finalize your choice by clicking **Create Project**, the working directory in RStudio will automatically be set to your project location. This means that any scripts, data files, or outputs you work on will be saved here by default, making them easier to find and reference later.

Inside your project directory, you’ll notice a file with an **.Rproj** extension, such as **MyProject.Rproj**. In the future, you can open this file to launch RStudio directly into this project, ensuring the working directory is already set. It’s also advisable to set up specific

folders within your project directory for different components like scripts, data, and figures. This keeps everything tidy and organized as your project expands.

2.8.3 Functions to import data in R

There are different functions to import data into R.

- For CSV (Comma Separated Values) files, the base R function `read.csv()` is commonly used. However, within the tidyverse package suite, the `readr` package provides the `read_csv()` function, which tends to be faster and more versatile. The `read_csv2()` function is designed for CSV files using semicolons ; as field separators and commas , as decimal points, compared to `read_csv()` which assumes commas as field separators.
- Excel files can be read using the `readxl` package, which provides the `read_excel()` function. This function can read both `.xls` and `.xlsx` files.
- For SPSS data files, you can use the `haven` package. This package contains the function `read_sav()` for `.sav` files.
- Stata data files, or `.dta` files, can also be read using the `haven` package with the `read_dta()` function.
- If you're working with R's native data format, `.RData` or `.RDS`, you can use the `load()` function for `.RData` files and the `readRDS()` function for `.RDS` files.
- Base R To import data in R, we will use the R base provides several functions to import data such as `read.csv()`. I personally prefer using the `readr` package which is part of the tidyverse. To read a csv file, you will need the `read_csv()` function.

All these functions enable you to read a file by specifying its path or its URL on the web. Here, for example, I import a CSV file containing data on candidates for the 2024 legislative elections in France.

```
library(tidyverse)

french_candidates <- read_csv("data/2024_french_candidates.csv")
```

```
Rows: 4009 Columns: 9
```

```
-- Column specification -----
```

```
Delimiter: ","
```

```
chr (7): code_department, departement, first_name, name, nuance, gender, pro...
```

```
dbl (2): exprimes_per, age
```

```
i Use `spec()` to retrieve the full column specification for this data.
```

```
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
glimpse(french_candidates)
```

```
Rows: 4,009
Columns: 9
$ code_department <chr> "01", "01", "01", "01", "01", "01", "01", "01", "01", ~
$ departement     <chr> "Ain", "Ain", "Ain", "Ain", "Ain", "Ain", "Ain", "Ain", "Ain"~
$ first_name      <chr> "Christophe", "Xavier", "Sébastien", "Vincent", "Éric"~
$ name            <chr> "MAÎTRE", "BRETON", "GUERAUD", "GUILLERMIN", "LAHY", "~
$ nuance          <chr> "RN", "LR", "UG", "ENS", "EXG", "DSV", "DSV", "RN", "D~
$ exprimes_per    <dbl> 39.37, 23.96, 23.45, 11.68, 0.69, 0.52, 0.33, 39.20, 2~
$ gender          <chr> "M.", "M.", "M.", "M.", "M.", "M.", "M.", "M.", "M.", ~
$ age             <dbl> 55, 62, 51, 48, 58, 36, 42, 35, 44, 37, 49, 69, 59, 43~
$ profession      <chr> "(22) - Commerçant et assimilé", "(33) - Cadre de la f~
```

2.8.4 Export data

Typically, we might want to save data to our disk after adding information, merging different datasets, and so on. This is useful for later reuse or to share the data with someone else. To achieve this, simply replace ‘read’ with ‘write’ in all the functions I’ve introduced previously. You’ll also need to specify the name of the R object containing your data and the path where you wish to export the data.

```
write_csv(french_candidates, "data/french_candidates2.csv")
```

2.9 Going further

2.9.1 Naming things

Note that I have written the names of objects with underscores. There are different conventions to write object names in R that you can discover [here](#). I personally use *snake case* which use lowercase letters and underscores to separate words.

2.9.2 Packages

Note above that I have used, the `::` operator. It lets you reference a specific function from a package without loading the whole package. This is handy when two packages have functions with the same name, ensuring clarity in your code. It’s also useful for one-off function uses, avoiding the need to load an entire package. This approach can make code clearer and sometimes faster by reducing loaded dependencies

2.9.3 Indexing

When we manipulate vectors, we often want to access specific elements of them, which we call **indexing**, which is performed by using square brackets `[]`. You can index either by position or by name.

When I write `leaders[3]`, I want the value of the third element of the `leaders` vector, this is indexing by position. But when I write `leaders[leaders == "Le Pen"]`, I index by name because I want the elements that have Le Pen as value.

```
leaders
```

```
[1] "Chassaigne" "Chatelain" "Panot"      "Vallaud"    "Ciotti"
[6] "Le Pen"
```

```
leaders[4] # Get the third element of the vector
```

```
[1] "Vallaud"
```

```
leaders[-3] # Get everything but the third element of the vector
```

```
[1] "Chassaigne" "Chatelain" "Vallaud"    "Ciotti"    "Le Pen"
```

```
leaders[c(1,4)] # Get the first and the fifth elements of the vector
```

```
[1] "Chassaigne" "Vallaud"
```

```
leaders[1:3] # Get elements from the first to the third
```

```
[1] "Chassaigne" "Chatelain" "Panot"
```

```
leaders[leaders == "Le Pen"] # Which has Le Pen as value
```

```
[1] "Le Pen"
```

```
leaders[leaders != "Le Pen"] # Which has not Le Pen as value
```

```
[1] "Chassaigne" "Chatelain" "Panot"      "Vallaud"    "Ciotti"
```

```
leaders[leaders %in% c("Le Pen", "Ciotti")]# Which has Le Pen or Ciotti
```

```
[1] "Ciotti" "Le Pen"
```

```
leaders[!leaders %in% c("Le Pen", "Ciotti")]# Which has neither Le Pen nor Ciotti
```

```
[1] "Chassaigne" "Chatelain" "Panot" "Vallaud"
```

2.9.4 More on vectors and functions

In R, functions often expect inputs of specific types. If you pass a character vector containing numeric numbers as strings to a function that expects a numeric vector, it may not behave as expected. As shown below, the function returns a NA which means Not available/applicable. When R encounters something it doesn't understand, it returns an error message with indications about the problem.

```
parl_seats_chr <- c("17", "38", "66", "72", "99", "36", "31", "47", "16", "126", "7", "22")
parl_seats_chr
```

```
[1] "17" "38" "66" "72" "99" "36" "31" "47" "16" "126" "7" "22"
```

```
mean(parl_seats_chr) # This returns NA
```

```
Warning in mean.default(parl_seats_chr): argument is not numeric or logical:
returning NA
```

```
[1] NA
```

Similarly, computing the sum of the `parl_seats_chr` vector will work as expected but trying to calculate the sum of our `leaders` character vector composed of leaders's names will not give a meaningful result.

```
sum(parl_seats_chr)
```

```
Error in sum(parl_seats_chr): invalid 'type' (character) of argument
```

```
sum(leaders) # This is an error
```

Error in sum(leaders): invalid 'type' (character) of argument

If you are not sure about the type of your vectors, you can check with the `class()` function that will give you the answer.

```
class(parl_seats)
```

```
[1] "numeric"
```

```
class(parl_seats_chr)
```

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

Sometimes, a vector has not the good type for the operation we want to perform. To check the type of a vector, you can use the family of `is.` functions such as `is.numeric()` and `is.character()` that return a boolean operator. In case the vector is not the right type for our purpose, we can try to coerce them with the family of `as.` functions such as `as.numeric()` and `as.character()`.

```
is.numeric(parl_seats_chr) # Check if numeric
```

```
[1] FALSE
```

```
parl_seats_num <- as.numeric(parl_seats_chr) # Coerce to numeric  
is.numeric(parl_seats_num) # Check again if numeric
```

```
[1] TRUE
```

```
mean(parl_seats_num) # Compute the mean
```

```
[1] 48.08333
```

Functions that you will find in R have been created by someone. You can also create your own functions in R. You usually start doing it when you are more advanced so do not worry if you find it hard, it is just for you to know that it is possible. Here I just create a simplified other function to calculate a mean in R.

```
# Create a function to compute the mean of a vector

compute_mean <- function(x) {
  # Compute the sum of the values in the vector and divide by the number of values (length)
  mean <- sum(x)/length(x)

  # Return the result
  return(mean)
}

compute_mean(parl_seats)
```

```
[1] NA
```

2.9.4.1 More on the tidyverse

The **tidyverse** isn't just a single package but rather a meta-package, meaning it bundles together several other packages, each with its own set of functions. For example, one of these bundled packages is **stringr**, which offers tools for manipulating character vectors. Since **stringr** is part of the **tidyverse**, if you've already loaded the **tidyverse**, there's no need to load **stringr** separately. With it, you can perform tasks like converting strings in a vector to uppercase or lowercase.

```
leaders
```

```
[1] "Chassaigne" "Chatelain"  "Panot"      "Vallaud"    "Ciotti"
[6] "Le Pen"
```

```
str_to_lower(leaders) # Change strings to lower class
```

```
[1] "chassaigne" "chatelain"  "panot"      "vallaud"    "ciotti"
[6] "le pen"
```

```
str_to_upper(leaders) # Change strings to upper class
```

```
[1] "CHASSAIGNE" "CHATELAIN"  "PANOT"      "VALLAUD"    "CIOTTI"
[6] "LE PEN"
```

```
str_detect(leaders, "C") # Detect if strings that contains a "C"
```

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

We can also combine characters vectors together with `str_c()`.

```
parties <- c("Communists", "Greens", "LFI", "Socialists", "Ciotti's party", "National Rally")  
str_c(leaders, " is the parliamentary leader of ", parties)
```

```
[1] "Chassaigne is the parliamentary leader of Communists"  
[2] "Chatelain is the parliamentary leader of Greens"  
[3] "Panot is the parliamentary leader of LFI"  
[4] "Vallaud is the parliamentary leader of Socialists"  
[5] "Ciotti is the parliamentary leader of Ciotti's party"  
[6] "Le Pen is the parliamentary leader of National Rally"
```


3 Getting help and dealing with errors

When learning a new language, mistakes are inevitable, and programming can be particularly frustrating due to the myriad of errors one can encounter. Unlike practicing a foreign language with speakers who may be forgiving of mistakes, R is exceptionally unforgiving. If your query is not structured correctly, it won't be understood at all. Therefore, it's crucial to find ways to pinpoint and understand the root causes of these errors. Overcoming errors is one of the most effective methods for learning R. Here are a few places where you can find the help you need:

3.1 R documentation

Every function or package in R comes with documentation provided by its creators. You can access this documentation directly from RStudio by placing a `?` before the function or package name in the console and executing the command. This will open the documentation in the help pane, where you'll find a description of the function, its various arguments, and some usage examples.

```
?readr::read_csv # Access the documentation of the read_csv function
```

3.2 Cheatsheets

Every package in the tidyverse (and some others) has a cheatsheet that provides information about its various functions. You can find these cheatsheets at this [link](#). For example, check out the `readr` [cheatsheet](#).

3.3 Online resources

First, remember that you're not alone in encountering errors in R; many others have faced similar issues before you. Often, they have sought help online. If R gives you an error you don't understand, it's likely that someone else has encountered the same issue and discussed it online. Start by checking if others have asked similar questions. Websites dedicated to R

programming frequently have solutions from experienced users. Additionally, Google can be a valuable resource. Whether you're trying to accomplish a task (e.g., “how to import a `.dta` file in R”) or resolve an error, searching the error message or your query can provide valuable insights.

Often, you'll find yourself on a site called [Stackoverflow](#), a community hub for users of various programming languages. You can often copy and paste the code you find there, but it's important to tweak and adapt it to fit your specific needs. Remember, in the world of programming, it's common for everyone to borrow and adapt code from others. Another helpful resource is the [Rstudio Community](#). Additionally, the [#rstats](#) hashtag on Twitter can provide insights and discussions from the R programming community.

3.4 AI is your friend

Lastly, we've entered the era of generative AI. When coding, Large Language Models, particularly [ChatGPT](#), can swiftly emerge as invaluable allies. By supplying ChatGPT with your errors or seeking guidance for specific coding tasks, you can obtain outstanding results. I strongly suggest incorporating it into your coding journey. However, exercise caution: ChatGPT might sometimes suggest non-existent functions or present inaccurate information. It isn't a magic bullet, but you can quickly assess the accuracy of its suggestions by testing the code in R. If the code doesn't work correctly, the information may be incorrect. For instance, check out this example [example](#) where I asked ChatGPT to explain how to import a Stata file in R.

3.5 Most common errors

Finally, some errors are really common and you will probably face them often. I provide you here a (non exhaustive) list of those to help you troubleshooting¹.

3.5.0.1 Syntax errors

Many errors beginners encounter in R stem from syntax issues: a slight coding mistake can lead RStudio to misunderstand your intentions. Common errors include typos in function names or forgetting symbols like `)`, `,`, or `"`. For example, if you missed a closing `"` when trying to subset the `Le Pen` string from the `leaders` vector: `leaders[leaders == "Le Pen]`, you'd likely see a `+` in the console. This indicates that R is awaiting further input to process your command.

¹I rely on many blogposts, [here](#), [here](#), [here](#), [here](#) and [here](#)

3.5.0.2 The “not found” errors

- Error: function 'x' not found : misspelling or package not loaded

```
Library(tidyverse)
```

```
Error in Library(tidyverse): could not find function "Library"
```

```
means(c(15,16,19))
```

```
Error in means(c(15, 16, 19)): could not find function "means"
```

```
read_html("https://labour.org.uk/category/latest/press-release/") # Read html code from a web
```

```
Error in read_html("https://labour.org.uk/category/latest/press-release/"): could not find f
```

Mistakes related to capitalization or misspelling are common. For instance, attempting to compute the mean of a number vector but mistakenly adding an “s” to the `mean()` function will lead to an error. In another scenario, you might aim to read a webpage’s HTML code for web scraping purposes. While the function might be correctly spelled, the error arises if the required `rvest` package isn’t loaded beforehand. When encountering such errors, ensure you’ve spelled functions correctly and loaded the necessary package (e.g., using `library(rvest)`).

- Error: object 'x' not found : typo, forgot to run the line or saving object

```
leaders <- c("Chassaigne", "Vallaud", "Chatelain", "Panot")
leader[1]
```

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

You might also want to look only at leaders from right-wing parties in the object `right_wing_leaders`. Here, the error happens because we did not save any object with this value yet.

```
right_wing_leaders
```

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

```
right_wing_leaders <- c("Waucquiez")
right_wing_leaders
```

```
[1] "Waucquiez"
```

Error in install.packages : object 'x' not found

```
install.packages(rvest)
```

Error in eval(expr, envir, enclos): object 'rvest' not found

Most of the time, you just forget the "" and you should write `install.packages("rvest")`. It might also be a typo in the package name (eg. you would have an error with `install.packag("Rvest")`).

- Error: 'x' does not exist in current working directory

```
readr::read_csv("thisdata.csv")
```

Error: 'thisdata.csv' does not exist in current working directory ('/Users/malo/Documents/te

This error typically arises when you've given an incorrect path, and R can't find your file. Use `getwd()` to check your current working directory and then adjust the file path as needed.

3.5.0.3 Inconsistent data types

We have seen already that R comes with different data types such as `logical` or `character`. Many functions takes as argument a vector of a specific type and will not work on other. Below an obvious example : if we try to compute the mean of a character vector, this will not work.

```
leaders
```

```
[1] "Chassaigne" "Vallaud"      "Chatelain"  "Panot"
```

```
class(leaders)
```

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

```
mean(leaders)
```

```
Warning in mean.default(leaders): argument is not numeric or logical: returning  
NA
```

```
[1] NA
```

Part II

Manipulating and describing data

Part III

Visualizing data

Part IV

Testing relationships

Part V

Regression analysis

Part VI

Multivariate analysis