

# **How to Use R for Data Science**

## **Lecture Notes**

Prof. Dr. Stephan Huber

May 12, 2025

# Table of contents

<b>I. Getting started with...</b>	<b>5</b>
<b>1. ...R</b>	<b>6</b>
1.1. Why R? . . . . .	6
1.2. How to learn R . . . . .	7
1.3. Learning resources . . . . .	9
1.4. What is a function in R? . . . . .	9
1.5. What are objects in R? . . . . .	11
1.6. What are R and RStudio? . . . . .	12
1.7. How to write and run code in R and RStudio . . . . .	14
1.8. How to install R, RStudio, and R packages . . . . .	15
1.9. What are R packages? . . . . .	16
1.9.1. Package installation . . . . .	16
1.9.2. Package loading . . . . .	17
1.9.3. Simplified package management with <code>p_load</code> . . . . .	17
1.10. Base R and the tidyverse universe . . . . .	19
1.11. Two key programming operators . . . . .	19
1.11.1. Assignment operator: “ <code>&lt;-</code> ” . . . . .	20
1.11.2. Pipe operator: “ <code> &gt;</code> ” . . . . .	20
1.12. Write your own function . . . . .	21
<b>2. ...writing code</b>	<b>22</b>
2.1. Bake a cake . . . . .	22
2.2. Elegant code . . . . .	24
2.3. Bake a cheese cake . . . . .	25
2.4. Comment what you do . . . . .	25
2.5. Bake 10 cakes . . . . .	26
2.6. Writing real code . . . . .	27
<b>3. ...writing R scripts</b>	<b>29</b>
3.1. The limitations of no-code applications . . . . .	29
3.2. R Scripts: Why they are useful . . . . .	30
3.3. Create, write, and run R scripts . . . . .	31
3.3.1. Create . . . . .	31
3.3.2. Write . . . . .	31
3.3.3. Run . . . . .	31
3.4. What to do at the header of each script . . . . .	32
<b>II. Basics of coding</b>	<b>34</b>
<b>4. Interactive introduction with swirl</b>	<b>35</b>
4.1. Set up <code>swirl</code> . . . . .	35
4.2. <code>swirl-it: huber-intro-1</code> . . . . .	36

## Table of contents

4.3. swirl-it: huber-intro-2 . . . . .	40
4.4. swirl-it: Data analytical basics . . . . .	45
4.5. swirl-it: The <code>tidyverse</code> package . . . . .	45
4.6. Other <code>swirl</code> modules . . . . .	45
<b>5. Kickstart</b>	<b>46</b>
5.1. Analysing the association of weight and the price of cars . . . . .	46
5.2. Accessing World Bank's <i>World Development Indicators</i> . . . . .	52
<b>6. Pitfalls</b>	<b>69</b>
6.1. No clue about the “working directory” . . . . .	69
6.2. No consistent directory structure . . . . .	69
6.3. Working manually outside R . . . . .	70
6.4. No active R Packages management . . . . .	70
6.5. Confusion between console and script . . . . .	70
6.6. Misunderstanding data types and formats . . . . .	71
6.7. Lack of knowledge about data identification . . . . .	71
6.8. Losing track of data due to excessive overwriting . . . . .	71
6.9. No documentation . . . . .	72
6.10. Ignoring error messages and warnings . . . . .	72
6.11. No attempt to identify the problem and troubleshoot . . . . .	73
6.12. Unstylish code . . . . .	76
<b>III. Do stuff</b>	<b>77</b>
<b>7. Manage data</b>	<b>78</b>
7.1. Import and generate data . . . . .	78
7.1.1. Assigning data to an object using the assignment operator <code>&lt;-</code> . . . . .	78
7.1.2. Vectors and matrices . . . . .	78
7.1.3. Open RData files . . . . .	80
7.1.4. Open datasets of packages . . . . .	80
7.1.5. Import data using public APIs . . . . .	80
7.1.6. Import various file formats . . . . .	82
7.1.7. Examples . . . . .	82
7.2. Data . . . . .	83
7.2.1. Data frames and tibbles . . . . .	83
7.2.2. Tidy data . . . . .	83
7.2.3. Data types . . . . .	85
7.3. Operators . . . . .	86
7.3.1. Algebraic operators . . . . .	86
7.3.2. The pipe operator: <code> &gt;</code> . . . . .	86
7.3.3. The <code>%in%</code> operator . . . . .	88
7.3.4. Extract operators . . . . .	88
7.3.5. Logical operators . . . . .	90
7.4. Data manipulation . . . . .	92
7.4.1. <code>dplyr</code> : A human readable grammar of data manipulation . . . . .	92
7.4.2. If statements . . . . .	97
7.4.3. Examining and cleaning data with the <code>janitor</code> package . . . . .	98
7.4.4. <code>tabyl()</code> - a better version of <code>table()</code> . . . . .	102
7.5. User-defined functions and conflicts . . . . .	103
7.6. Example: How to explore a dataset . . . . .	105

*Table of contents*

<b>8. Visualize data</b>	<b>109</b>
<b>9. Collection of exercises</b>	<b>111</b>
9.1. Import data with <code>c()</code> . . . . .	111
9.2. Filter and select observations . . . . .	112
9.3. Base R, <code>%in%</code> operator, and the pipe <code> &gt;</code> . . . . .	112
9.4. Generate and drop variables . . . . .	118
9.5. Subsetting . . . . .	121
9.6. Weighting, data management and regression . . . . .	125
9.7. Consumer prices over time . . . . .	128
9.8. Load the Stata dataset “auto” using R . . . . .	138
9.9. DatasauRus . . . . .	143
9.10. Convergence . . . . .	151
9.11. Unemployment and GDP in Germany and France . . . . .	162
9.12. Import data and write a report . . . . .	173
9.13. Explain the weight of students . . . . .	176
9.14. Calories and weight . . . . .	186
9.15. Bundesliga . . . . .	194
9.16. Okun’s Law . . . . .	219
9.17. Names and duplicates . . . . .	230
9.18. Zipf’s law . . . . .	238
<b>IV. Publish</b>	<b>256</b>
<b>10. Quarto</b>	<b>257</b>
10.1. Why Markdown and Quarto? . . . . .	257
10.1.1. No-code vs. code-based writing applications . . . . .	257
10.1.2. Typical (mis)usage of WYSIWYG applications . . . . .	258
10.1.3. Advantages of Quarto for writing text . . . . .	258
10.2. Markdown . . . . .	259
10.3. Quarto . . . . .	259
10.3.1. Introduction . . . . .	259
10.3.2. Create an APA compliant manuscript using Quarto . . . . .	261
10.4. R Markdown . . . . .	262
<b>11. Collaborating with Git and GitHub</b>	<b>265</b>
11.1. Introduction . . . . .	265
11.2. Install Git . . . . .	266
11.3. Using Git from the terminal . . . . .	267
11.3.1. Configuring Git . . . . .	267
11.3.2. Initializing a Repository . . . . .	268
11.3.3. Staging Changes . . . . .	268
11.3.4. Committing Changes . . . . .	269
11.3.5. Pushing Changes . . . . .	269
11.3.6. Undo changes . . . . .	269
11.4. Using Git from RStudio . . . . .	270
11.4.1. Set up Git in RStudio . . . . .	270
11.4.2. Connecting RStudio Projects with GitHub repositories . . . . .	270
11.5. Make a contribution using Git and GitHub . . . . .	271

*Table of contents*

<b>12. Create and host a website</b>	<b>274</b>
12.1. Creating a website with Quarto . . . . .	274
12.2. Hosting the website on GitHub . . . . .	274
<b>References</b>	<b>276</b>
<b>Appendices</b>	<b>278</b>
<b>A. Navigating the file system</b>	<b>278</b>
A.1. The file system . . . . .	278
A.2. Working directory . . . . .	278
A.3. Navigating the file system using the R console . . . . .	279
A.4. R Studio projects . . . . .	280
A.5. Why do the Windows paths use the back-slash? . . . . .	280
<b>B. Operators</b>	<b>281</b>
B.1. Assignment: . . . . .	281
B.2. Arithmetic: . . . . .	281
B.3. Relational: . . . . .	281
B.4. Logical: . . . . .	281
B.5. Others: . . . . .	281
<b>C. Popular functions</b>	<b>282</b>
C.1. Help . . . . .	282
C.2. Package management . . . . .	282
C.3. General . . . . .	282
C.4. Tools . . . . .	282
C.5. Data import . . . . .	282
C.6. Inspect data . . . . .	283
C.7. Graphics . . . . .	283
C.8. Data management . . . . .	283
C.9. <code>dplyr</code> functions . . . . .	284
C.10. Data analysis . . . . .	284
C.11. Statistical functions . . . . .	284
<b>D. Helpful shortcuts</b>	<b>286</b>

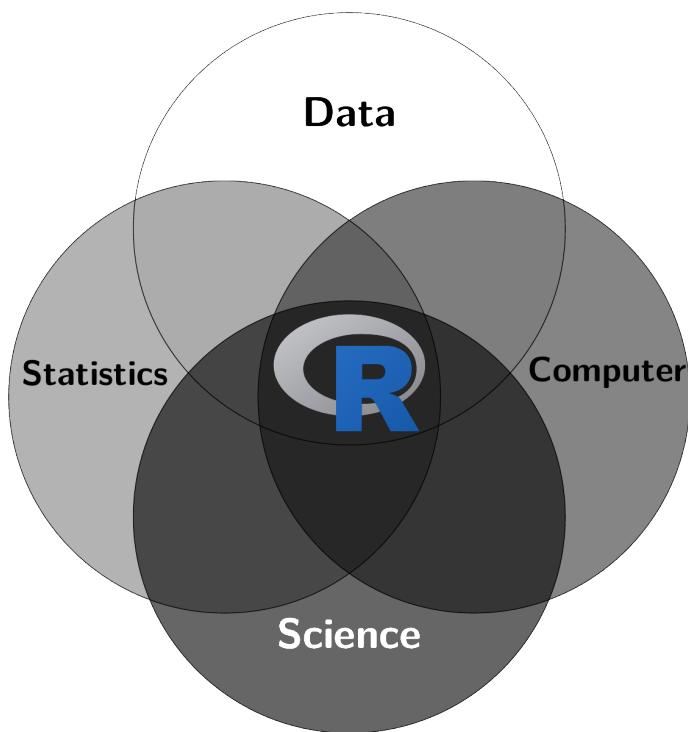
# List of figures

1.	Prof. Dr. Stephan Huber . . . . .	3
1.1.	Play around with code . . . . .	8
1.2.	A collection of textbooks . . . . .	10
1.3.	The R Documentation of <code>q()</code> . . . . .	12
1.4.	Analogy of difference between R and RStudio . . . . .	13
1.5.	Icons of R versus RStudio on your computer . . . . .	13
1.6.	A sketch of RStudio interface to R . . . . .	13
1.7.	One plus four in a R script . . . . .	15
1.8.	Set up R in three steps . . . . .	15
1.9.	The tidyverse universe . . . . .	19
6.1.	Googlling the error message . . . . .	73
7.1.	The logo of the packages <code>readr</code> , <code>haven</code> , and <code>readxl</code> . . . . .	82
7.2.	The logos of the <code>tidyverse</code> and <code>tibble</code> packages . . . . .	83
7.3.	Features of a tidy dataset: variables are columns, observations are rows, and values are cells . . . . .	84
7.4.	The logo of the <code>dplyr</code> package . . . . .	92
8.1.	Pie charts are problematic . . . . .	109
9.1.	The logo of the <code>DatasauRus</code> package . . . . .	143
9.2.	Weight vs. Calories . . . . .	188
9.3.	Ranking history: 1. FC Kaiserslautern . . . . .	196
9.4.	Ranking history: 1. FC Köln . . . . .	196
10.1.	Example of an R Markdown file . . . . .	262
10.2.	R Markdown Cheatsheet from Posit . . . . .	262
10.3.	Xie et al. [2020]: R Markdown Cookbook . . . . .	263
10.4.	Xie et al. [2018]: R Markdown: The Definitive Guide . . . . .	263
11.1.	The FINAL.doc problem . . . . .	265
11.2.	GitHub is big . . . . .	265
11.3.	Memorizing six git commands . . . . .	266
11.4.	Three git commands you really need . . . . .	267
11.5.	Copy the https URL of your repo . . . . .	269
11.6.	Fork the repo . . . . .	271

# List of tables

6.1. Typical folder structure . . . . .	69
6.2. Most frequent errors . . . . .	72
7.1. Basic algebraic operators . . . . .	86
7.2. Logical operators . . . . .	90
9.1. Covid cases and deaths till August 2022 . . . . .	111
9.2. Data . . . . .	126
11.1. Most important git commands . . . . .	267
11.2. Most common bash commands . . . . .	268
D.1. Different OS, different keys . . . . .	286
D.2. Helpful shortcuts . . . . .	286

# Preface



## About R

The programming language R enables you to handle, visualize, and analyze data. It is compatible with various operating systems (Windows, Mac, Linux) and can do a lot of things better compared to other programs like Python, Stata, Eviews, SPSS, SAS, and Excel. R is open source, extensively utilized, and there are abundant resources available for learning it. These notes are just my five cents.

## About the cover of the notes

Data science is a buzzword that combines different fields of knowledge such as computer science, software engineering, informatics, database management, statistics, econometrics, business intelligence, and mathematics. However, there is no universally accepted definition of it and I think it is not important to define it precisely. [Kelleher and Tierney \[2018, p. 97\]](#) wrote “Data science is best understood as a partnership between a data scientist and a computer.” So data science is about embracing the power of computers for scientific, commercial or social purposes. Of course, empirical models and statistics play a role in gaining meaningful insights. The graphic on the cover page may illustrate that R combines four important fields, that are, data, science, computer, and statistics.

## About the notes

 A PDF version of these notes is available [here](#).

Please note that while the PDF contains the same content, it has not been optimized for PDF format. Therefore, some parts may not appear as intended.

- These notes aims to support my lecture at the HS Fresenius but are incomplete and no substitute for taking actively part in class.
- I hope you find this book helpful. Any feedback is both welcome and appreciated.
- This is work in progress so please check for updates regularly.
- These notes offer a curated collection of explanations, exercises, and tips to facilitate learning R without causing unnecessary frustration. However, these notes don't aim to rival comprehensive textbooks such as [Wickham and Grolemund \[2023\]](#).
- These notes are published under the [Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License](#). This means it can be reused, remixed, retained, revised and redistributed as long as appropriate credit is given to the authors. If you remix, or modify the original version of this open textbook, you must redistribute all versions of this open textbook under the same license. This script draws from the work of [Navarro \[2020\]](#), [Muschelli and Jaffe \[2022\]](#), [Thulin \[2021\]](#), and [Ismay and Kim \[2022\]](#)



which is also published under the same license.

- I host the notes in a [GitHub repo](#).

 To reap the best benefits from studying,

I recommend to copy all the code that is shown in the book into a R script and try to run it on your PC. That is the best way to learn, understand, and create your own notes that may guide you later on. Whenever you see interesting code somewhere, try to run it on your PC. Moreover, I recommend the exercises of the book, they are challenging sometimes but to really understand code you need to run code yourself.

## Structure of these notes

Chapter	Explanations
...R	Learn the basics everyone should know about R and RStudio, including how to install them.
...writing code	Learn the basics of writing code.
...writing R scripts	Learn how to use R scripts and their benefits.
Interactive introduction using swirl	A hands-on tutorial on how to use the swirl package. This section is optional.
Kickstart	A quick start guide for beginners on how to dive into R, showcasing some of its capabilities.
Pitfalls	Discover common mistakes beginners often make and how to avoid them to save time on troubleshooting.
Manage data	Learn how to manipulate data in R.
Visualize data	A quick guide on where to find resources to learn about creating graphical visualizations in R.
Collection of exercises	A set of exercises to practice R programming skills.

Chapter	Explanations
Appendix	A set of useful stuff that will help you to navigate through your file system, find the right operator and function, or to learn some useful shortcuts.

## About the author

### Contact:

Prof. Dr. Stephan Huber  
Hochschule Fresenius für Wirtschaft & Medien GmbH  
Im MediaPark 4c  
50670 Cologne  
Office: 4e OG-3  
Telefon: +49 221 973199-523  
Mail: [stephan.huber@hs-fresenius.de](mailto:stephan.huber@hs-fresenius.de)  
Private homepage: [www.hubchev.github.io](http://www.hubchev.github.io)  
Github: <https://github.com/hubchev>

Figure 1.: Prof. Dr. Stephan Huber



I am a Professor of *International Economics and Data Science* at HS Fresenius, holding a Diploma in Economics from the University of Regensburg and a Doctoral Degree (summa cum laude) from the University of Trier. I completed postgraduate studies at the Interdisciplinary Graduate Center of Excellence at the Institute for Labor Law and Industrial Relations in the European Union (IAAEU) in Trier. Prior to my current position, I worked as a research assistant to Prof. Dr. Dr. h.c. Joachim Möller at the University of Regensburg, a post-doc at the Leibniz Institute for East and Southeast European Studies (IOS) in Regensburg, and a freelancer at Charles University in Prague.

Throughout my career, I have also worked as a lecturer at various institutions, including the TU Munich, the University of Regensburg, Saarland University, and the Universities of Applied Sciences in Frankfurt and Augsburg. Additionally, I have had the opportunity to teach abroad for the University of Cordoba in Spain, the University of Perugia in Italy, and the Petra Christian University in Surabaya, Indonesia. My published work can be found in international journals such as the Canadian Journal of Economics and the Stata Journal. For more information on my work, please visit my private homepage at [hubchev.github.io](http://hubchev.github.io).

## Preface

I was always fascinated by data and statistics. For example, in 1992 I could name all soccer players in Germany’s first division including how many goals they scored. Later, in 2003 I joined the introductory statistics course of [Daniel Rösch](#). I learned among others that probabilities often play a role when analyzing data. I continued my data science journey with [Harry Haupt’s Introductory Econometrics](#) course, where I studied the infamous Jeffrey M. [Wooldridge \[2002\]](#) textbook. It got me hooked and so I took all the courses [Rolf Tschernig](#) offered at his chair of Econometrics, where I became a tutor at the University of Regensburg and a research assistant of [Joachim Möller](#). Despite everything we did had to do with how to make sense out of data, we never actually used the term *data science* which is also absent in the more 850 pages long textbook by [Wooldridge \[2002\]](#). The book also remains silent about *machine learning* or *artificial intelligence*. These terms became popular only after I graduated. The *Harvard Business Review* article by [Davenport and Patil \[2012\]](#) who claimed that data scientist is “The Sexiest Job of the 21st Century” may have boosted the popularity.

The term “data scientist” has become remarkably popular, and many people are eager to adopt this title. Although I am a professor of *data science*, my professional identity is more like that of an applied, empirically-oriented international economist. My hesitation to adopt the title “data scientist” also stems from the deep respect I have developed through my interactions with econometricians and statisticians. Considering their in-depth expertise, I feel like a passionate amateur.

Ultimately, I poke around in data to find something interesting. Much like my ten-year-old younger self who analyzed soccer statistics to gain a deeper understanding of the sport. The only thing that has changed since then is that I know more promising methods and can efficiently use tools for data processing and data analysis.

**Part I.**

**Getting started with...**

# 1. ...R

## Learning Objectives

- Understand the reasons for choosing R as a programming language.
- Learn effective strategies and resources for mastering R programming.
- Learn the basic principles of R.
- Distinguish between R and RStudio.
- Demonstrate how to write and execute code in R and RStudio.
- Learn how to install R, RStudio, and R packages.
- Explore options to use R and RStudio without installing them on your machine, such as through cloud-based platforms.

## 1.1. Why R?

R is an open-source programming language that allows to analyse and manipulate data, create state-of-the-art graphics, and many more. It supports larger data sets, reads any type of data, and runs on multiple platforms (Windows, Mac, Linux) and CPU architectures (x86\_64, arm64). R makes it easier to automate tasks, organize projects, ensure reproducibility, and find and fix errors, and anyone can contribute packages to improve its functionality. Moreover, the following points are worth to emphasize:

- **R is an artist!** Check out:
  - [The R Graph Gallery](#)
  - [R CHARTS by R CODER](#)
- **R is an employment insurance!** Programming is a core skill in research, economics, and business. If you can write code, you have plenty of opportunities to earn a decent salary. R is one of the most widely used programming languages in the world today. It is used in almost every industry such as finance, banking, medicine or manufacturing. R is used for portfolio management, risk analytics in finance and banking industries. Even if you need to learn a new programming language later, knowing R makes it much easier to pick up another one.
- **R uses the computer and computers are great!** Doing statistics on a computer is faster, easier and more powerful than doing it by hand. Computers are an extension to your brain and can do repetitive tasks better and faster without making logical errors. The only reason to do statistical calculations with pencil and paper is for learning purposes.
- **Low-code and no-code applications such as Excel are limited!** Using spreadsheets software like Microsoft Excel for research can be problematic. It's easy to lose track of operations, making the process difficult to oversee and document. Command-line programs are maybe not as easy to learn but offer a more straightforward approach that allows the results to be replicated easily.

- **R is open source!** Proprietary software expansive, support can only be provided by the copyright owner which means the software expires and you can't do anything against it. Moreover, security issues cannot be checked as the source code is not available, and possibilities for customization are limited. R is yours and everybody can contribute to its success.
- **R is big!** When you download and install R, you get some basic packages, that contain functions that allow you to do already a lot of things. Beyond that, you can write your own packages or install user-written packages that extend your possibilities. With over 20,684 packages on the CRAN repository and many more available on GitHub and other platforms, R's extensive library supports a wide variety of data science tasks. Its widespread use and open-source availability have cemented R as a standard tool in data science and ensured that there are multiple approaches to most data handling processes. These can be easily adopted.



### R has weaknesses

For newcomers to programming, the learning curve is rather flat at the beginning. One reason is that R tools are spread across many packages, which can overwhelm beginners. There is no centralized support and the helpful and active online community have different backgrounds. It can be difficult for beginners to find the right solution as there are often many different ways to tackle the same problem. Moreover, R can be slower than languages like Python, MATLAB, C/C++ or Java.

## 1.2. How to learn R

There are many different approaches to learning R. It pretty much depends on your preferences, needs, goals, prerequisites and limitations. It is up to you to search and find a suitable way to achieve your learning goals. While I hope you find my notes helpful, I additionally provide in section [Section 1.3](#) a list of other resources that are worth considering. To start with, I recommend my swirl courses that provide an interactive learning environment, see [Chapter 4](#).

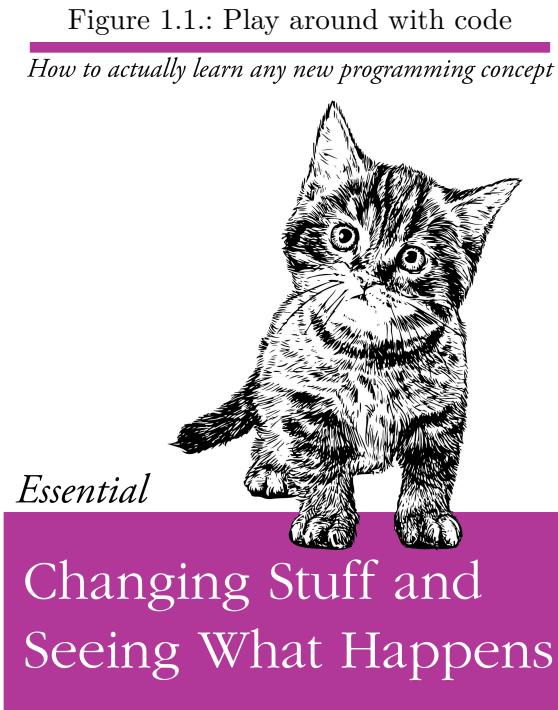


### Make your hands dirty!

Learning a programming language can, like learning a foreign language, be daunting and frustrating. However, if you put in the effort and are not afraid to make mistakes, anybody can learn it. You don't have to be a nerd. To have a guide next to you can help and speed up your progress significantly. The key is taking action and getting involved. I mean, do write code. Try to copy the code that you read here and elsewhere. Explore what the code does on your machine. Don't be afraid to make errors. Your PC will not explode. In this paper, most of the code is written in a manner that allows you to effortlessly copy and reproduce the output on your PC. Take advantage of this opportunity and go for it! Hands-on practice is far more enjoyable than merely reading through the material.

Here are some comments that may help you to learn efficiently:

- **Computers need clear and precise instructions to work:** They can't handle mistakes or unclear directions. They are actually sort of stupid as they do not have an intuition. They just take you literally. Even small errors like a missing comma or an unclosed bracket can cause your code to not work. Computers do exactly what you tell them, no more and no less.



O RLY?

@ThePracticalDev

*Source: DEV Community on GitHub*

### i Computers take you literally

Let me illustrate what I mean: Suppose you send your grandfather the following message:

“Let’s eat grandpa.”

He will probably understand that you’re inviting him to dinner. However, if you sent the same message to a computer, it would interpret the sentence literally due to the missing comma:

“Let’s eat, grandpa.”

The comma makes all the difference in clarifying that you’re speaking to your grandpa, not about eating him! Similarly, in programming, an incorrectly placed comma can break your code or change the meaning of your code.

- **Copy, paste, and tweak:** While learning code from scratch is sometimes essential, you can speed up your work by modifying code that already exists. I call this the “*copy, paste, and tweak*” approach. While this is not the only way to learn code, it gets a job done quick, and it is fun, see Figure 1.1.
- **Have a purpose when coding:** Rather than learning to code for its own sake, it is more fun and you’ll probably learn faster when you have a goal in mind. Try to analyze data that you are interested in. Another good exercise is replicating a research paper.
- **Practice is key:** The best method to improving your coding skills is through lots of practice. Consequently, these notes give you plenty of exercises.

- **Use ChatGPT:** The usage of supporting tools is not forbidden. ChatGPT can help you to understand code and brainstorm solutions. However, it's important to know that ChatGPT might suggest complex methods when there are shorter and more elegant solutions available. Absolute beginners might find ChatGPT's solutions overwhelming and have difficulties to tweak the proposed sketch of a solution. So, use it thoughtfully.

### 1.3. Learning resources

Thousands of freely available books and resources exist. [bookdown.org](#) and the [Big Book of R](#) are two vast collections of links to R books that might verify my claim.

In RStudio you find in the right side at the bottom a panel that is called *Help*. There you find a lot of links, manuals, and references that offer you tons of resources to learn R for free including: [education.rstudio.com](#) and [Links for Getting Help with R](#). At the top right of RStudio you find a panel called tutorial. Here you can install the `learnr` package that offers some nice interactive tutorials.

Since you may feel overwhelmed by the number of resources, I would like to highlight some books:

1. [Wickham and Grolemund \[2023\]: R for Data Science: Import, Tidy, Transform, Visualize, and Model Data](#) is the most popular source to learn R. It focuses on introducing the tidyverse package and is freely available online.
2. [Healy \[2018\]: Data Visualization: A Practical Introduction](#) is a hands-on introduction to the principles and practice of looking at and presenting data using R and `ggplot`.
3. [Irizarry \[2022\]: Introduction to Data Science: Data Analysis and Prediction Algorithms With R](#) is a complete, up to date, and applied introduction.
4. [Venables et al. \[2022\] An Introduction to R: Notes on R: A Programming Environment for Data Analysis and Graphics](#) is a manual from the R Core Development Team that shows how to use R without having to install and load additional packages.
5. [Neth \[2023\]: Data Science for Psychologists](#) is a comprehensive introduction to R and data science for non experts of both programming and data science. It uses a variety of data types and includes many examples and exercises.
6. [Kabacoff \[2024\]: Modern Data Visualization with R](#) teaches how to create graphs from scratch providing a lot of examples that you can copy, paste and tweak.

Some other sources that are worth mentioning are these:

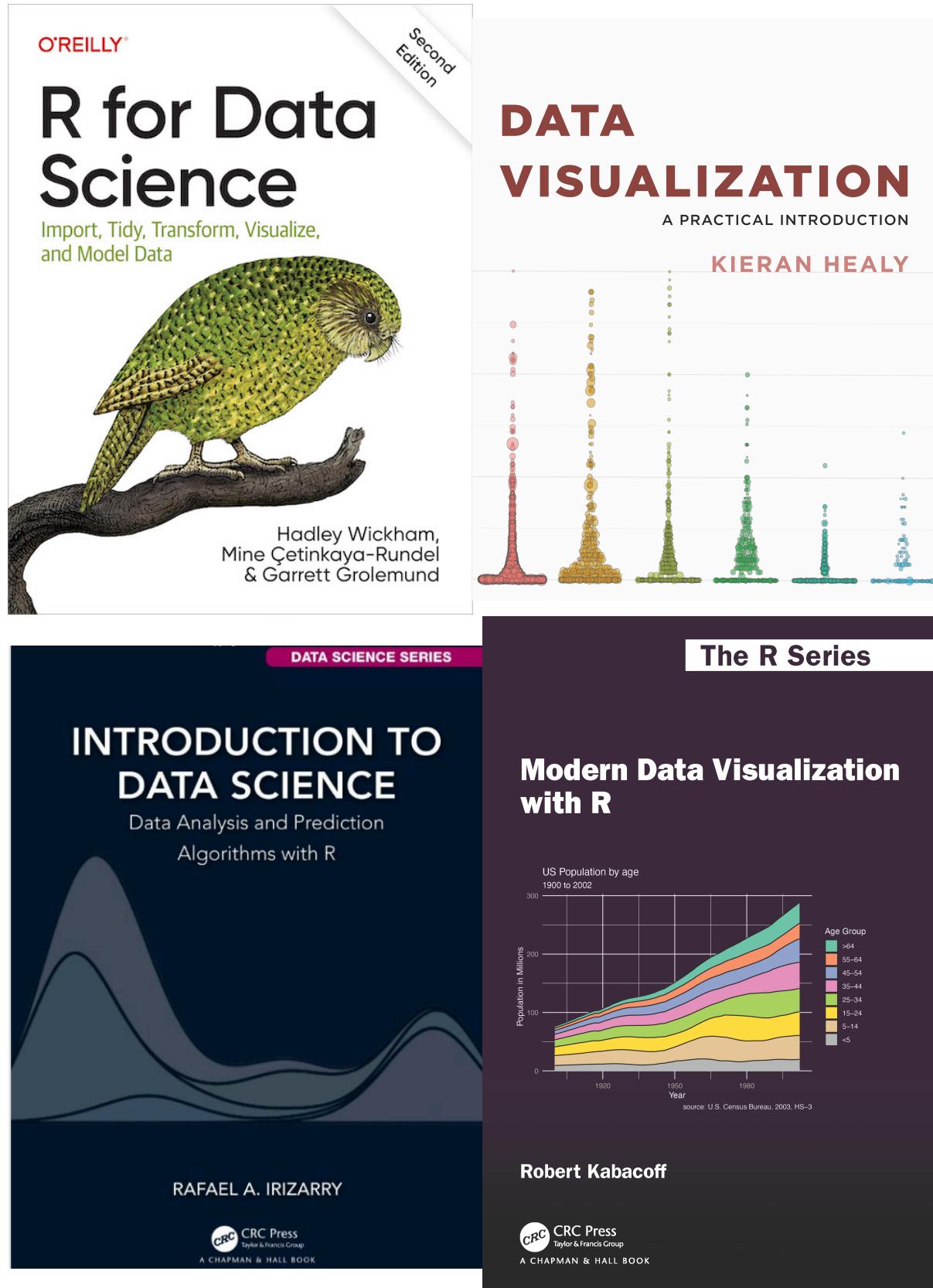
- The search engine [www.rseek.org](#) is R specific and often better than [www.google.com](#) as it only searches for content that has to do with the programming language R.
- On [rdocumentation.org](#) you can find the complete documentation of all R packages.
- Many find these [cheatsheets](#) helpful.

### 1.4. What is a function in R?

R is a *functional programming language*. If you want R to do something, you need to use a function. Or, in the words of [Chambers \[2017\]](#), p. 4]:

“Everything that happens is a function call.”

Figure 1.2.: A collection of textbooks



For example, when you like to exit R, you do it with the function `q()`:

```
> q()
Save workspace image? [y/n/c] :
```

If you want to specify what exactly you want R to do for you, you need to refer to the arguments of a function. For example, if you don't want to be asked interactively what you want to do with your workspace (this is the place where you store all your objects, see section [Section 1.5](#)), you can do this with an argument that is part of the `q()` function:

```
> q(save = "no")
```

To learn more about a function, you can access its documentation by typing a question mark followed by the function name into the Console:

```
?q()
```

Unfortunately, the documentation can sometimes be a bit confusing for beginners in applied contexts. However, the documentation for all functions is structured similarly, typically featuring several key sections:

- **Description:** A brief overview of what the function does.
- **Usage:** How to use the function, including the function name and its arguments.
- **Arguments:** Detailed descriptions of each argument the function accepts, including what types of values are expected.
- **Details:** Additional details about the function's behavior and any important notes.
- **Examples:** Practical examples demonstrating how to use the function in various contexts.

Understanding these sections can significantly enhance your ability to navigate and utilize R.

An excerpt of the *R Documentation* for the function `q()` is shown in [Figure 1.3](#). Here, we observe that the function has three arguments that you can manipulate. If you do not specify any of these arguments explicitly, we see that *by default*, R sets the three arguments as shown.

## 1.5. What are objects in R?

R is an object oriented programming language. That means,

“everything that exists in R is an object” [[Chambers, 2017](#), p. 4].

Objects are the fundamental units that are used to store information. Objects can be a variety of data types, including vectors, matrices, data frames, lists, functions. Moreover, you can store empirical results, tables, figures and many more in form of so-called objects. All objects are shown in the *workspace* which is shown in the *Environment* panel.

In R, you can show the content of the workspace with `ls()`. The function `rm()` allows to remove objects and with `rm(list=ls())` you clear all objects from the workspace.

Figure 1.3.: The R Documentation of `q()`

**Description**

The function `quit` or its alias `q` terminate the current R session.

**Usage**

```
quit(save = "default", status = 0, runLast = TRUE)
      q(save = "default", status = 0, runLast = TRUE)
```

**Arguments**

- `save` a character string indicating whether the environment (workspace) should be saved, one of "no", "yes", "ask" or "default".
- `status` the (numerical) error status to be returned to the operating system, where relevant. Conventionally 0 indicates successful completion.
- `runLast` should `.Last()` be executed?

**Details**

`save` must be one of "no", "yes", "ask" or "default". In the first case the workspace is not saved, in the second it is saved and in the third the user is prompted and can also decide *not* to quit. The default is to ask in interactive use but may be overridden by command-line arguments (which must be supplied in non-interactive use).

## 1.6. What are R and RStudio?

While R has a command line interface, there are multiple third-party graphical user interfaces available that improve the user experience a lot. The most successful graphical user interface or integrated development environment (IDE) is RStudio. Throughout this book, I will assume that you are using R via RStudio. First time users often confuse the two. At its simplest, R is like a car's engine while RStudio is like a car's dashboard as illustrated in Figure Figure 1.4.

More precisely, R is a functional programming language that runs computations, while RStudio is an *integrated development environment (IDE)* that provides an interface by adding many convenient features and tools. So just as the way of having access to a speedometer, rearview mirrors, and a navigation system makes driving much easier, using RStudio's interface makes using R much easier as well.

Much as we don't drive a car by interacting directly with the engine but rather by interacting with elements on the car's dashboard, we won't be using R directly but rather we will use RStudio's interface. After you install R and RStudio on your computer, you'll have two new *programs* (also called *applications*) you can open. We'll always work in RStudio and not in the R application. Figure Figure 1.5 shows what icon you should be clicking on your computer.

After you open RStudio, you should see something similar to Figure Figure 1.6 where three or four panels dividing the screen.

1. The *Environment* panel, where a list of all objects is shown.
2. The *Files*, *Plots* and *Help* panel, allow you to manage files, preview plots, and find help for different functions of R.

Figure 1.4.: Analogy of difference between R and RStudio

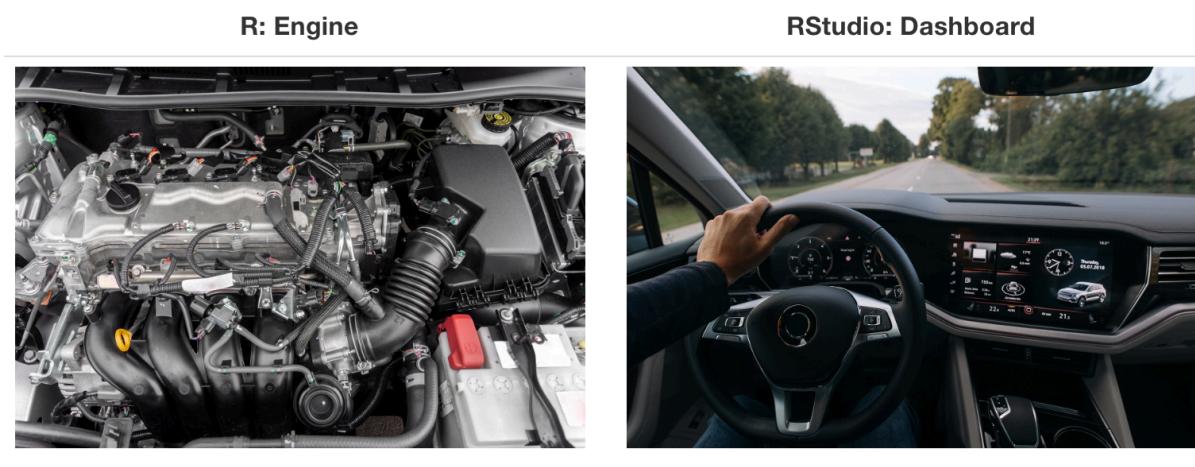


Figure 1.5.: Icons of R versus RStudio on your computer

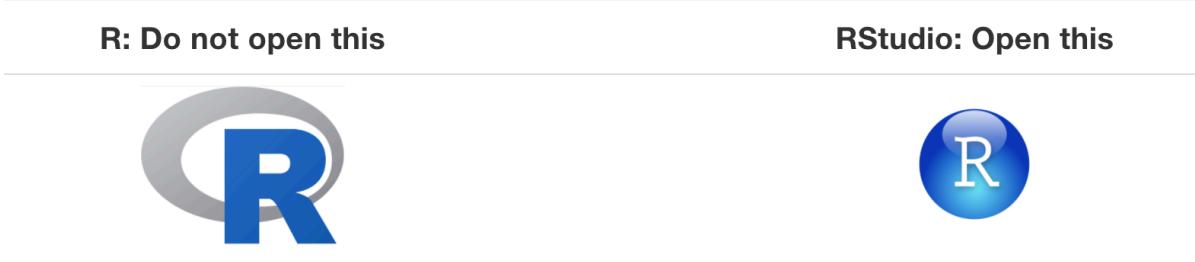
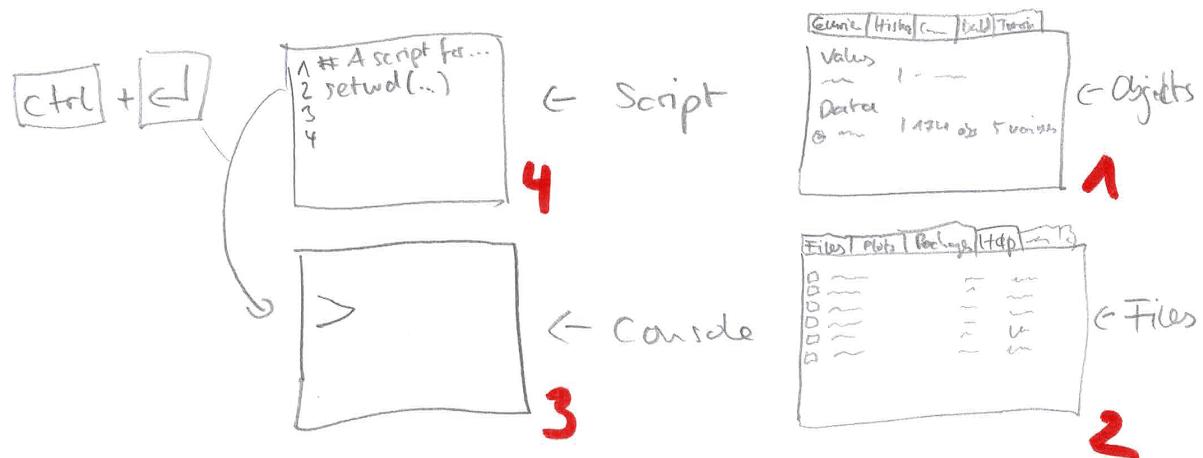


Figure 1.6.: A sketch of RStudio interface to R



3. The *Console* panel, used for running code.
4. The *Script* panel, used for writing code.

If you don't have panel number 4,

open it by opening an existing R-script or creating a new one. You can create a new on by clicking *Ctrl+Shift+N* (alternatively, you can use the menu: File→New File→R Script).

The *Console* panel will contain R's startup message, which shows information about which version of R you're running. My startup message at the time of writing was as follows:

```
R version 4.3.3 (2024-02-29) -- "Angel Food Cake"
Copyright (C) 2024 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (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.

Natural language support but running in an English locale

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.
```

You can resize the panels as you like, either by clicking and dragging their borders or using the minimise/maximise buttons in the upper right corner of each panel. Clicking *Ctrl++* and *Ctrl+-* allows to make the fonts larger or smaller.

## 1.7. How to write and run code in R and RStudio

In the Console you can type in code and push Enter to run the line of code. For example, you can calculate:

**1+4**

[1] 5

While working in the Console is possible, we usually work in RStudio using so-called *scripts*. These scripts are plain text files with the file extension ".R". Scripts are discussed in detail in Chapter 3. To create a script, go to the *File* menu, select *New File* and then choose *R Script*.

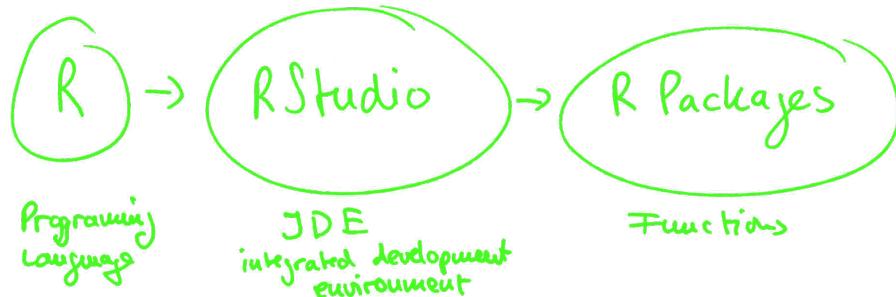
With the key shortcut **Ctrl+Enter** for Windows and Linux user or by **Cmd+Enter** for MacOs users (or by clicking **Run**) you can run a line of a script, that means you send one line of code to the Console. See Figure 1.7 how this looks like in RStudio.

Figure 1.7.: One plus four in a R script

The screenshot shows the RStudio interface. At the top is a menu bar with File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help. Below the menu is a toolbar with various icons. A script editor window titled "Untitled1" is open, containing the code "1 1+4" on line 1 and "2" on line 2. Below the script editor is a status bar showing "2:1 (Top Level) R Script". At the bottom is a tab bar with "Console" (selected), "Terminal", "Markers", and "Background Jobs". The console tab shows the command "> 1+4" followed by the output "[1] 5".

## 1.8. How to install R, RStudio, and R packages

Figure 1.8.: Set up R in three steps



As shown in Figure 1.8, setting up R on your personal computer (Windows, Mac, Linux) is a three step process: You will first need to download and install R. After that has been successful you can download and install RStudio. Please note that it is important that you install R first and then install RStudio. As a third but optional step you can install R packages.

### 1. Do this firstly: Download and install R [here](#).

- If you are a Windows user: Click on “Download R for Windows”, then click on “base”, then click on the Download link.
- If you are macOS user: Click on “Download R for (Mac) OS X”, then under “Latest release:” click on R-X.X.X.pkg, where R-X.X.X is the version number. For example, the latest version of R as of March 29, 2024 was R-4.3.3.
- If you are a Linux user: Click on “Download R for Linux” and choose your distribution for more information on installing R for your setup.

### 2. Do this secondly: Download and install RStudio [here](#).

- Scroll down to “Installers for Supported Platforms” near the bottom of the page.
- Click on the download link corresponding to your computer’s operating system.

### 3. Do this thirdly: Install R packages. This step is optionally as you can install R packages at any time. However, it may be a good idea to install frequently used packages in one take

because the installation of some packages can be time consuming. Therefore, I recommend to read Section 1.9 and follow the instructions therein.

### 💡 How to use R and RStudio without installation

If you don't want to install R on your PC or you don't have admin rights to do so or if you want to run R on your tablet (IPad or Chromebook) or even your smartphone, you can use RStudio online doing *cloud computing* on <https://posit.cloud>. Posit Cloud (formerly RStudio Cloud) is a cloud-based solution that allows anyone to use RStudio online and navigate it through your web browser. It is free for individuals with some restrictions and limited capacities.

## 1.9. What are R packages?

A package is a collection of functions, data sets and other R objects that are all grouped together under a common name. More than 20,000 packages are available at the official repository (CRAN). CRAN is a network of ftp and web servers around the world that store identical, up-to-date, versions of code and documentation for R, see: <https://cran.r-project.org>].

However, before we get started, there's a critical distinction that you need to understand, which is the difference between having a package **installed** on your computer, and having a package **loaded** in R. When you install R on your computer only a small number of packages come bundled with the basic R installation. The installed packages are on your computer. The critical thing to remember is that just because something is on your computer doesn't mean R can use it. In order for R to be able to *use* one of your installed packages, that package must also be *loaded*. Generally, when you open up R, only a few of these packages (about 7 or 8) are actually loaded.

### ℹ️ Package management

1. A package must be installed before it can be loaded.
2. A package must be loaded before it can be used.

We only need to install a package once on our computer. However, to use the package, we need to load it every time we start a new R environment or R Studio, respectively.

### 1.9.1. Package installation

To install an R package you can use the GUI of R Studio or the command line. In R Studio you can click on the *Packages* tab, then on the *Install* button, then you must search for a package and click *Install*. An alternative way to install a package is by typing

```
install.packages("package_name")
```

in the console pane of RStudio and pressing Return/Enter on your keyboard. Note you must include the quotation marks around the name of the package.

If you want to update a previously installed package to a newer version, you need to re-install it by repeating the earlier steps or you use `update.packages()`. To uninstall packages you can use `remove.packages()`.

 How to speed up the installation of packages

The installation of packages can take some time. However, if your CPU has many cores, you can speed up the process a lot using the argument `Ncpus` like this `update.packages(ask = F, Ncpus = 4L)`. This option allows you to adjust the number of parallel processes R can use on your PC. So, if you have a CPU with many cores you can increase that number. A tutorial on how to set the number of cores used by R permanently can be found [here](#).

### 1.9.2. Package loading

Recall that after you've installed a package, you need to *load it*. We do this by using the `library()` command. For example, to load the `ggplot2` package, run the following code in the console pane. What do we mean by "run the following code"? Either type or copy-and-paste the following code into the console pane and then hit the Enter key.

```
library("ggplot2")
```

If after running the earlier code, a blinking cursor returns next to the > "prompt" sign, it means you were successful and the `ggplot2` package is now loaded and ready to use. If, however, you get a red "error message" that reads

```
Error in library(ggplot2) : there is no package called 'ggplot2'
```

It means that you didn't successfully install it. If you get this error message, go back to section Section 1.9.1 on R package installation and make sure to install the `ggplot2` package before proceeding.

One very common mistake new R users make when wanting to use particular packages is they forget to *load* them first by using the `library()` command we just saw. Remember: *you have to load each package you want to use every time you start RStudio*. If you don't first *load* a package, but attempt to use one of its features, you'll see an error message similar to:

```
Error: could not find function
```

R is informing you that you are attempting to use a function from a package that has not yet been loaded. Forgetting to load packages is a common mistake made by new users, and it can be a bit frustrating to get used to at first. However, with practice, it will become second nature for you. Unloading packages can be done with `detach(package:ggplot2, unload=TRUE)`.

### 1.9.3. Simplified package management with p\_load

I recommend to install and load packages using the `p_load()` function of the `pacman` package. It is superior because

- it only installs a package if it is has not been installed yet,
- it loads the package, and
- does not require quotes nor the `c()`function.

For example, instead of the traditional approach:

```
install.packages(
  c("tidyverse", "janitor", "haven", "readxl")
)
library(
  c("tidyverse", "janitor", "haven", "readxl")
)
```

You can streamline the process as follows:

```
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, janitor, haven, readxl)
```

The line `if (!require(pacman)) install.packages("pacman")` ensures the installation of the `pacman` package, which is necessary for using the `p_load` function.

Before you load packages in a script, I recommend to *unload* all other packages with

```
pacman::p_unload(all)
```

to avoid conflicts of functions (see Section 7.5).

### Tip 1: Install everything now

Throughout the lecture notes and in the exercises, I will use different packages. The installation can be time consuming and hence I recommend to install all packages by running the following lines of code in the Console. This takes some minutes depending on your PC and your internet connection. However, after installing all these packages you have all packages that are used in my exercises, my lecture notes *How to Use R for Data Science*, and the book *R for Data Science (2e)* by Wickham and Grolemund [2023].

```
if (!require(pacman)) install.packages("pacman")
pacman::p_load(
  arrow, babynames, car, curl, datasauRus, devtools, dplyr,
  duckdb, devtools, expss, gapminder, ggplot2, ggrepel, ggridges,
  ggpubr, ggstats, ggthemes, haven, HH, janitor, kableExtra, knitr,
  Lahman, labelled, likert, magick, maps, MASS, nycflights13,
  openxlsx, palmerpenguins, papaja, plm, psych,
  remotes, rempsyc, repurrrsive, rstatix, skimr, sjlabelled,
  sjmisc, sjPlot, stargazer, texreg, tidyR,
  tidyverse, tinylabels, usethis, WDI, wbstats, writexl
)
```

In addition to these packages, I recommend to install a package that I created to offer you some tutorials and functions. I host this package on my GitHub account and you can install it as follows:

```
devtools::install_github("hubchev/hubchev")
```

## 1.10. Base R and the tidyverse universe

Upon successfully installing R, you gain access to functions that are part of *Base R*. This includes standard packages automatically installed and loaded with each R session, such as `stats`, `utils`, and `graphics`, providing a broad spectrum of functionalities for statistical analysis and graphical capabilities [see [Venables et al., 2022](#)]. However, the syntax in *Base R* can become complex and less intuitive for users. Consequently, many individuals, including Hadley Wickham, the Chief Data Scientist at *Posit* (formerly RStudio), and his team, have developed an alternative suite of packages known as the `tidyverse`. These packages share a common philosophy and syntax, emphasizing readability and ease of use. We will heavily utilize the `tidyverse` in the following sections.

Figure 1.9.: The tidyverse universe



The R package `tidyverse` (see Figure 1.9) is a comprehensive collection of R packages including popular packages such as `ggplot2`, `dplyr`, `tidyr`, `readr`, `purrr`, `tibble`, `stringr`, and `forcats`, which together offer extensive capabilities for data modeling, transformation, and visualization.

How to do data science with `tidyverse` is the subject of multiple books and tutorials. In particular, the popular book *R for Data Science* by [Wickham and Grolemund \[2023\]](#) is all about the `tidyverse` universe. Thus, I highly recommend reading sections [Workflow: basics](#)), [Data transformation](#), and [Data tidying](#). Additionally, explore [www.tidyverse.org](http://www.tidyverse.org) for more resources, and consider completing the `tidyverse` module in my `swirl` package, `swirl-it`, as detailed in section Chapter 4.

To install and load `tidyverse` run the following lines of code:

```
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse)
```

**Exercise 1.1.** Set up R, RStudio, and R packages  
Open [this interactive tutorial](#) and work through it.

## 1.11. Two key programming operators

To instruct R to perform a task, we use function calls. When we want R to utilize data, we refer to that data within an object. This leads to two important questions:

1. How do we create objects?
2. How can we instruct R to execute multiple steps in sequence?

### 1.11.1. Assignment operator: “<-”

The assignment operator “<-” is used to store data in an object or overwrite an existing object. For example, we can calculate the square root of 5 using the `sqrt()` function and assign the result to an object named `square_root_of_five` as follows:

```
square_root_of_five <- sqrt(5)
```

Now, if you call the object, R will return the result:

```
square_root_of_five
```

```
[1] 2.236068
```

The assignment operator is also explained in Section [7.1.1](#).

### 1.11.2. Pipe operator: “|>”

There are different ways to chain function calls in R. The base R package allows you to nest one function within another. For example, to calculate the sum of the square root of 5 and the square root of 9 using the `sum()` function, you can write:

```
sum(sqrt(5), sqrt(9))
```

```
[1] 5.236068
```

In this case, we sum the square roots of 5 and 9, with the two functions nested as arguments within `sum()`. If you want to round the result to two decimal places, you can use the `round(x, digits = 2)` function like this:

```
round(sum(sqrt(5), sqrt(9)), digits = 2)
```

```
[1] 5.24
```

This is another example of function nesting.

Alternatively, you can use the pipe operator, which can be represented as “|>” in base R or as “%>%” if you load the `magrittr` package. The pipe operator passes the output of one function to serve as the input for the next. Here’s how you can perform the same calculation using the pipe operator:

```
c(5, 9) |>
  sqrt() |>
  sum() |>
  round(digits = 2)
```

```
[1] 5.24
```

This can be interpreted step by step as follows:

- `c(5, 9)` : We combine the values 5 and 9 into a vector using the `c()` function ...AND THEN...
- `sqrt()`: We calculate the square root of each value in the vector ...AND THEN...
- `sum()`: We sum the square roots ...AND THEN...
- `round(digits = 2)`: We round the result to two decimal places.

As you can see, the pipe operator allows us to read the code as “and then”. This method of sequentially executing tasks has several advantages: it mimics how humans typically approach problems and makes the code easier to read and understand. Consequently, we will frequently utilize this operator throughout the book. The pipe operator is also explained in Section [7.3.2](#).

## 1.12. Write your own function

Defining your own functions in R is straightforward. For example, if you frequently need to perform the calculation described in Section [1.11](#), you can create a custom function like this:

```
process_numbers <- function(num1, num2) {
  num1 |>
  sqrt() |>
  sum() |>
  round(digits = 2)
}
```

The function `process_numbers()` has two arguments, that are two numbers, as input, performs the calculation, and returns the result. Now, you can use this function to process any two numbers with ease. For example, for the numbers 5 and 9, the function call is:

```
process_numbers(5, 9)
```

```
[1] 5.24
```

And for the numbers 4 and 9, it is:

```
process_numbers(4, 9)
```

```
[1] 5
```

This approach not only simplifies your code but also ensures consistency when performing the same operation multiple times.

How to define user-defined function is explained in greater detail in Section [7.5](#).

## 2. ...writing code

When you talk to a computer, you use a programming language. Computers can handle certain tasks for you, like doing math or creating graphs. They don't complain and work fast. They just need clear instructions. Giving clear instructions to a machine, however, isn't easy. Unlike humans, computers lack intuition and cannot adapt to the context of your commands; they interpret instructions literally. Check out [this video](#). It'll show you that good communication among human beings isn't about explaining things perfectly, but about using your gut feeling and common sense effectively.

If you want a computer to handle computationally intensive and repetitive tasks, you must learn to "speak" in a way that the computer understands. Never assume the computer knows what you mean. The following sections will emphasize this: Be precise and be specific.

### 2.1. Bake a cake

As a teacher, I often contemplate how I can teach students to communicate with a computer. In this section, I will attempt to translate a cake baking recipe into a programming language, using R. I hope you find this both amusing and insightful. Let's start by examining a simplified cake recipe:

#### Instructions to bake a cake:

1. Buy an oven.
2. Buy ingredients (flour, sugar, butter, eggs, soda).
3. Buy tools.
4. Clean everything..
  
5. Heat up the oven.
6. Prepare the tools (springform, bowl, mixer).
7. Weight all ingredients.
8. Take a bow and all the weighted ingredients and put everything in a bowl.
9. Take the mixer and mix all ingredients in the bowl for 3 minutes.
10. Put all the mixed ingredients in a springform pan.
11. Take the springform pan with the mixed ingredients and put it in the oven.
12. Bake for 30 minutes, take the springform pan it out of the oven, and turn off the oven.
13. Clean the kitchen and the tools.

Although this recipe is simplified, it illustrates a process you might be familiar with. Now, let's assume a computer is tasked with baking a cake. How would we explain the necessary steps to the computer using the R programming language?

## 2. ...writing code

In R, a functional programming language, we understand that everything that happens is a function call, and everything that exists is an object. Therefore, we must translate all actions into functions and all items into objects.

Here's what the translated recipe could look like in R:

```
buy(oven, springform, bowl, mixer, flour, sugar, butter, eggs, soda)
clean(oven, springform, bowl, mixer)
turn_on(oven)
prepare(springform, bowl, mixer)
weigh(flour, sugar, butter, eggs, soda)
dough <- bowl |>
  put(flour, sugar, butter, eggs, soda) |>
  action(tool = mixer, time = 3)

dough_springform <- springform |>
  put(dough) |>

dough_oven <- oven |>
  put(dough_springform) |>
  action(tool = oven, time = 30) |>
  pull()

turn_of(oven)
clean(oven, springform, bowl, mixer)
```

Understanding the translation of a recipe into code becomes clearer when we familiarize ourselves with two key programming operators:

1. The “`<-`” is known as the *assignment operator*. It saves or stores data into a new object. It might be helpful to think of it as saying, “I create the object `<name of object>` and store therein”
2. The “`|>`” is known as the *pipe operator*. It passes the output of one action to serve as the input for the next. Think of it as saying “and then.”

For example, the following lines:

```
dough <- bowl |>
  put(flour, sugar, butter, eggs, soda) |>
  action(tool = mixer, time = 3)
```

can be interpreted as:

```
I create the object `dough` and I store therein the bowl, and then
I put flour, sugar, butter, eggs, and soda to it, and then
I take action with the mixer for 3 minutes
```

In the preceding functions, you'll notice objects separated by commas and parameters like `tool = mixer, time = 3`. These parameters define the behavior of the function. When there's nothing within the brackets, as in `pull()`, the input is merely the output of the preceding pipe operator.

Even though R is no good as a cook and the recipe is missing some steps, this analogy helps to illustrate how programming languages work: they allow us to instruct the computer in a sequential way. Next, I will showcase why coding is appealing.

## 2.2. Elegant code

Let's make our code more *elegant*, that is, easy to read, understand, and modify. For example, while it is equivalent to write everything in one line

```
dough <- bowl |> put(flour, sugar, butter, eggs, soda) |> action(tool = mixer, time = 3)
```

or spread out over three lines,

```
dough <- bowl |>
  put(flour, sugar, butter, eggs, soda) |>
  action(tool = mixer, time = 3)
```

it is easier for the human eye to read the text in spread out form.

### 💡 Style in Writing Code

Writing code involves certain conventions, often referred to as a *coding style*. Although not strictly necessary, a consistent style can significantly enhance clarity and prevent common pitfalls. Numerous style guides aim to standardize coding practices. For example, you might find [The tidyverse style guide](#) by Hadley Wickham particularly helpful in adopting a harmonious coding style in R.

By using the assignment operator `<-`, we can create two objects: `ingredients` and `tools`. These objects are used multiple times throughout the process.

Here is an improved version of the script:

```
ingredients <- c(flour, sugar, butter, eggs, soda)
tools <- c(oven, springform, bowl, mixer)

buy(tools, ingredients)

clean(tools)
turn_on(oven)
prepare(tools)
weight(ingredients)
dough <- bowl |>
  put(ingredients) |>
  action(tool = mixer, time = 3)

dough_springform <- springform |>
  put(dough)

dough_oven <- oven |>
  put(dough_springform) |>
```

## 2. ...writing code

```
action(tool = oven, time = 30) |>
  pull()

turn_of(oven)
clean(, tools)
```

This version refines the process, making the code more streamlined and easier to follow.

### 2.3. Bake a cheese cake

Now, let's assume you want to bake another cake, this time with chocolate and banana, but without eggs. Moreover, you need to bake it for 45 minutes. We can easily adapt the code snippet from above to accommodate the ingredients for this new recipe:

```
ingredients <- c(flour, sugar, butter, soda, banana, chocolate)
tools <- c(oven, springform, bowl, mixer)

buy(tools, ingredients)

clean(tools)
turn_on(oven)
prepare(tools)
weight(ingredients)
dough <- bowl |>
  put(ingredients) |>
  action(tool = mixer, time = 3)

dough_springform <- springform |>
  put(dough)

dough_oven <- oven |>
  put(dough_springform) |>
  action(tool = oven, time = 45) |>
  pull()

turn_of(oven)
clean(kitchen, tools)
```

### 2.4. Comment what you do

Sometimes code can be difficult to understand for humans. It is therefore helpful to add comments to clarify what the individual code sections are supposed to do. In R, comments can be added with a leading hashtag, #.

```
# Decide on tools and ingredients
ingredients <- c(flour, sugar, butter, soda, banana, chocolate)
tools <- c(oven, springform, bowl, mixer)
```

```

# Go shopping
buy(tools, ingredients)

# Prepare the kitchen, tools, and ingredients
clean(tools)
turn_on(oven)
prepare(tools)
weight(ingredients)

# Make the dough
dough <- bowl |>
  put(ingredients) |>
  action(tool = mixer, time = 3)
dough_springform <- springform |>
  put(dough) |>

# bake the cake
dough_oven <- oven |>
  put(dough_springform) |>
  action(tool = oven, time = 45) |>
  pull()

# Clean up
turn_of(oven)
clean(kitchen, tools)

```

## 2.5. Bake 10 cakes

As a computer can reproduce a cake within seconds (I mean, not really, just in my little fun exercise here), we now have the opportunity to experiment with several versions of the cake by varying the baking time from 35 to 45 minutes. Here's how the corresponding code might look:

```

ingredients <- c(flour, sugar, butter, soda, banana, chocolate)
tools <- c(springform, bowl, mixer)

buy(tools, ingredients)

clean(tools)
turn_on(oven)
prepare(tools)
weight(ingredients)
dough <- bowl |>
  put(ingredients) |>
  action(tool = mixer, time = 3)

dough_springform <- springform |>
  put(dough)

```

## 2. ...writing code

```
for (timing in 35:44) {  
  dough_oven <- oven |>  
  put(dough_springform) |>  
  action(tool = oven, time = timing) |>  
  pull()  
  assign(paste("dough_oven_min_", timing, sep = ""), dough_oven)  
}  
  
turn_of(oven)  
clean(tools)
```

You can see a loop with some new and tweaked lines:

```
for (timing in 35:44) {  
  dough_oven <- oven |>  
  put(dough_springform) |>  
  action(tool = oven, time = timing) |>  
  pull()  
  assign(paste("dough_oven_min_", timing, sep = ""), dough_oven)  
}
```

These lines sequentially execute the following actions:

```
Let the object timing be 35, make a cake, and save it in the object `dough_oven_min_35` then  
let the object timing be 36, make a cake, and save it in the object `dough_oven_min_36` then  
let the object timing be 37, make a cake, and save it in the object `dough_oven_min_37` then  
let the object timing be 38, make a cake, and save it in the object `dough_oven_min_38` then  
let the object timing be 39, make a cake, and save it in the object `dough_oven_min_39` then  
let the object timing be 40, make a cake, and save it in the object `dough_oven_min_40` then  
let the object timing be 41, make a cake, and save it in the object `dough_oven_min_41` then  
let the object timing be 42, make a cake, and save it in the object `dough_oven_min_42` then  
let the object timing be 43, make a cake, and save it in the object `dough_oven_min_43` then  
let the object timing be 44, make a cake, and save it in the object `dough_oven_min_44` then
```

After all, we have ten cakes. This shows how we can harness the processing power of a computer. Computers are excellent at performing everyday, repetitive tasks so that we can automate processes and perform procedures effortlessly over and over again.

## 2.6. Writing real code

Of course, computers can't *bake a cake*. The R programming language can do none of the above. Nevertheless, there are analogies to the programming language R. Let me present a few lines of code and explain these lines of code to you, and you will see that the similarities are striking.

Copy that code chunk, paste it into a R script and run it.

## 2. ...writing code

```
# This script demonstrates a typical data analysis workflow in R
# -----
#
# Install and load required libraries
if (!require(pacman)) install.packages("pacman")
pacman::p_unload(all)
pacman::p_load(tidyverse,haven, janitor)

# Set the working directory to a project-specific folder
setwd("~/Documents")

# Clear the current environment of any objects
rm(list = ls())

# Load data from a Stata file available online
auto <- read_dta("http://www.stata-press.com/data/r18/auto.dta")

# Display basic information about the dataset
ncol(auto) # Number of columns
nrow(auto) # Number of rows
dim(auto) # Dimensions of the dataset
names(auto) # Names of variables
head(auto) # First few rows
tail(auto) # Last few rows
summary(auto) # Summary statistics for each column
glimpse(auto) # Compact display of the structure of the dataset
print(auto, n = Inf) # Print all rows of the dataset

# Check for duplicate entries based on the 'make' variable
auto |>
  get_dupes(make)

# Create and display a scatter plot of car price versus weight
plot_weight_price <- ggplot(auto, aes(x = weight, y = price)) +
  geom_point()
plot_weight_price

# Save the plot to a file
ggsave("plot_weight_price.png", plot = plot_weight_price, dpi = 300)
```

## 3. ...writing R scripts

### 3.1. The limitations of no-code applications

No-Code Applications (NCA) such as [Microsoft Excel](#), [RapidMiner](#), [KNIME](#), [DataRobot](#), [Tableau](#), [Microsoft Power BI](#), and [Google AutoML](#) are popular for good reasons. They enable the application of advanced empirical methods with no or minimal programming effort. Their intuitive graphical user interfaces come with pre-built templates and drag-and-drop functionality which helps to get things done quick, without having to study the program documentation for hours. Despite their apparent ease of use and efficiency, these platforms come with several disadvantages compared to traditional ways of working with computer by scripting and coding. Understanding these weaknesses helps to see why professional researchers, especially those actively publishing in academic journals, tend to rely on scripting languages like R and Python. These programming languages offer full control, are customizable and extendable, offer extensive opportunities for automation and reproducibility, and are often better suited for demanding data science tasks.

#### Disadvantages of No-Code Applications (NCA)

- NCA lack flexibility and are limited in their adaptability.
- NCA often have problems scaling with increasing data volumes or user requirements.
- NCA are often closed systems that make it difficult to integrate other systems or applications.
- NCA often bind a company to a specific ecosystem. This can lead to dependencies and vendor lock-in.
- NCA can obscure the underlying logic of how applications work, which can make troubleshooting more difficult.
- NCA abstracts the coding process and prevents users from understanding fundamental concepts that could be beneficial to their professional growth and ability to tackle more complex problems.

In summary, while low-code and no-code platforms offer quick deployment and ease of use, they can lack the depth, flexibility, and control provided by traditional scripting. Researchers and businesses must consider these trade-offs, especially when planning for long-term scalability, complex customizations, or in-depth integrations.

#### A hypothetical but realistic example with Excel

Suppose you work with a spreadsheet software like Excel. You import a CSV file using the implemented import tool, you save the converted file. You notice that Excel has messed up the dates during the import, so you spend a few minutes cleaning that manually. Then, you visualize the data and you save the visualizations in various tabs. Maybe, you can spot some outliers and you document in footnote that these outliers are the result of some issues with the raw data. As these errors cannot be solved, you delete the observations and variables that contain a significant amount of errors. Finally, you use filters to calculate

### 3. ...writing R scripts

some summary statistics. All your results are written in a new tab. You're convinced that you've done a great job. However, you send the file to your supervisor and you ask him for her opinion.

- She probably asks you what you have done to the data. How can you efficiently and completely communicate that?
- She comes back to you some time later and asks you to do the same analysis with an updated version of the data. How can you exactly redo the analysis and how long do you think it will take you to complete the job?
- She finds an error in your work or she has an idea to improve your analysis by making a few adjustments. Can you implement the adjustments easily, or do you have to redo everything from scratch?
- She sends you back the file with the comment that she has worked out some things in the file and now everything should be fine. How do you know what she has done to the data?

To say it in the words of [Stephenson \[2023\]](#), sec. 5.1]:

*“Spreadsheets are a nightmare for quality control and reproducibility, and you should always think twice before using one. Spreadsheets will always be a handy way to manipulate tabular datasets, and you’ll probably find them useful for data collection and quick back-of-the-envelope calculations, but they’re often more trouble than they’re worth.”*

## 3.2. R Scripts: Why they are useful

I have already discussed in Section 1.7 that you can run code either directly by typing your code into the console of R or by writing a script and then sending the code with **Ctrl+Enter** or with the **Run** button to the console of R. Typing functions into the console to run code may seem simple, but this interactive style has limitations:

- Typing commands one at a time can be cumbersome and time-consuming.
- It's hard to save your work effectively.
- Going back to the beginning when you make a mistake is annoying.
- You can't leave notes for yourself.
- Reusing and adapting analyses can be difficult.
- It's hard to do anything except the basics.
- Sharing your work with others can be challenging.

That's where having a transcript of all the code, which can be re-run and edited at any time, becomes useful. An R script is just plain text that is interpreted as code or as a comment if the text follows a hastag **#**. A script comes with important advantages.

Scripts...

- ... provide a record of everything you did during your data analysis.
- ... can easily be edited and re-run.
- ... allow you to leave notes for yourself.
- ... make it easy to reuse and adapt analyses.
- ... allow you to do more complex analyses.
- ... make it easy to share your work with others.

## 3.3. Create, write, and run R scripts

### 3.3.1. Create

**i** 3 equivalent ways to create a script

1. Use the menu: **File > New File > R Script**
2. Use the keyboard shortcut: **Ctrl+Shift+N** (Windows/Linux) or **Cmd+Shift+N** (Mac) or
3. Type the following in the console:

```
file.create("hello.R")
```

In the first two ways, a new R script window will open which can be edited and should be saved either by clicking on the **File** menu and selecting **Save**, clicking the disk icon, or by using the shortcut **Ctrl+S** (Windows/Linux) or **Cmd+S** (Mac). If you go for the third way, you need to open it manually.

### 3.3.2. Write

Regardless of your preferred way of generating a script, we can now start writing our first script:

```
x <- "hello world"
print(x)
```

Then save the script using the menus (**File > Save**) as **hello.R**.

The above lines of code do the following:

- With the assignment operator `<-` we create an object that stores the words “hello world” in an object entitled `x`. In Section 7.1.1 the assignment operator is further explained.
- With the third input we print the content of the object `x`.

### 3.3.3. Run

So how do we run the script? Assuming that the `hello.R` file has been saved to your working directory, then you can run the script using the following command:

```
source( "hello.R" )
```

Suppose you saved the script in a sub-folder called *scripts* of your working directory, then you need to run the script using the following command:

```
source("./scripts/hello.R")
```

### 3. ...writing R scripts

Just note that the dot, ., means the current folder. Instead of using the `source` function, you can click on the `source` button in Rstudio.

With the character # you can write a comment in a script and R will simply ignore everything that follows in that line onwards.

#### Exercise 3.1. Run a script and round numbers

Please copy and paste the following lines of code into an R script, run it on your computer, and try to understand how it works.

```
# Create a vector that contains the sales data
sales_by_month <- c(0, 100, 200, 50, 3, 4, 8, 0, 0, 0, 0)
sales_by_month
sales_by_month[2]
sales_by_month[4]
february_sales <- sales_by_month[2]
february_sales
sales_by_month[5] <- 25 # added May sales data
sales_by_month
# Do I have 12 month?
length( x = sales_by_month )
# Assume each unit costs 7 Euro, then the revenue is
price <- 7
revenue <- sales_by_month*price
revenue
# To get statistics for daily revenue we define the number of days:
days_per_month <- c(31, 28, 31, 30, 31, 30, 31, 31, 30, 31, 30, 31)
# Calculate the daily revenue
revenue_per_day <- revenue/days_per_month
revenue_per_day
# round number
round(revenue_per_day)
```

Use the “?” to search for the documentation of all functions used. In particular, do you understand how the function `round()` works? What arguments does the function contain? How can you manipulate the pre-defined arguments. For example, can you calculate the rounded revenue per day with two or four digits? Try it out!

```
?round()
```

#### Solution

```
round(revenue_per_day, digits = 4)
```

## 3.4. What to do at the header of each script

At the beginning of each script, ensure that all required packages are loaded correctly. Use the `pacman` package, which provides the `p_load()` function to load and, if necessary, install packages, and the `p_unload(all)` function to unload all packages. Additionally, set your work-

### 3. ...writing R scripts

ing directory with `setwd()` and clear all objects from the environment with `rm(list = ls())`. This ensures that everything in the environment after sourcing the script originates from the script itself. Below is the code that I use at the beginning of all my scripts. I recommend you do the same.

💡 Tip 5: Start your script with

```
if (!require(pacman)) install.packages("pacman")
pacman::p_unload(all)
pacman::p_load(tidyverse, janitor)
setwd("~/your-directory/")
rm(list = ls())
```

## **Part II.**

# **Basics of coding**

## 4. Interactive introduction with `swirl`

This section is designed to kickstart your journey into data science with R through the R package `swirl` that offers an interactive learning platform. `swirl` teaches you R programming and data science interactively, at your own pace, and right in the R console! You get immediate feedback on your progress. If you are new to R, have no fear. `swirl` will walk you through each of the steps required to employ Rstudio and R for your purpose.

For those seeking additional or alternative resources beyond `swirl`, exploring other introductory textbooks and resources on R is highly recommended. Please consider the resources I discuss in Section 1.3. One notable example is [Irizarry \[2022\]](#) who provides a comprehensive and conservative approach to understanding R.

### 4.1. Set up `swirl`

To install `swirl` and my learning modules, please follow my instructions precisely!

Open Rstudio and type in the console the following:

```
install.packages("swirl")
library("swirl")
install_course_github("hubchev", "swirl-it")
swirl()
```

The above four lines of code do the following:

- Install the `swirl` package, ensuring it's available for use in R.
- Load the `swirl` package, making its functions accessible.
- Install my `swirl` course that is hosted on GitHub, making its functions accessible.
- By entering `swirl` into the Console (located at the bottom-left in RStudio) and pressing the Enter key, you initiate `swirl`. This begins your interactive learning experience with the package.

 Tip 3: If the course has failed to install,

you can try to download the file `swirl-it.swc` from [github.com/hubchev/swirl-it](https://github.com/hubchev/swirl-it) and install the course with loading the `swirl` package and typing `install_course()` into the console.

After initiating the `swirl` environment, follow the instructions displayed in the Console. Specifically, select the *swirl-it* course and the *huber-intro-1* learning module to begin. You can exit `swirl` at any moment by typing `bye()` into the Console or pressing the *Esc* key on your keyboard.

## 4.2. swirl-it: huber-intro-1

**i** Click to see the full content of the module

Welcome to this swirl course. If you find any errors or if you have suggestions for improvement, please let me know via [stephan.huber@hs-fresenius.de](mailto:stephan.huber@hs-fresenius.de).

The RStudio interface consists of several windows. You can change the size of the windows by dragging the grey bars between the windows. We'll go through the most important windows now.

Bottom left is the Console window (also called command window/line). Here you can type commands after the > prompt and R will then execute your command. This is the most important window, because this is where R actually does stuff.

Top left is the Editor window (also called script window). Here collections of commands (scripts) can be edited and saved. When you do not get this window, you can open it with 'File' > 'New' > 'R script'.

Just typing a command in the editor window is not enough, it has to be send to the Console before R executes the command. If you want to run a line from the script window (or the whole script), you can click 'Run' or press 'CTRL+ENTER' to send it to the command window.

The shortcut to send the current line to the console and run it there is \_\_\_\_\_.

- a) CTRL+SHIFT
- b) CTRL+ENTER
- c) CTRL+SPACE
- d) SHIFT+ENTER

*Hint: You find all shortcuts in the menu at Tools > Keyboard Shortcuts Help or click ALT+SHIFT+K. If you are a Mac user, your shortcut is 'Cmd+Return' instead of 'SHIFT+ENTER'. To move on type skip().*

### Solution

answer: b

Top right is the environment window (a.k.a workspace). Here you can see which data R has in its memory. You can view and edit the values by clicking on them.

Bottom right is the plots / packages / help window. Here you can view plots, install and load packages or use the help function.

The first thing you should do whenever you start Rstudio is to check if you are happy with your working directory. That directory is the folder on your computer in which you are currently working. That means, when you ask R to open a certain file, it will look in the working directory for this file, and when you tell R to save a data file or figure, it will save it in the working directory.

You can check your working directory with the function `getwd()`. So let's do that. Type in the command window `getwd()` .

`getwd()`

```
[1] "/home/sthu/Dropbox/hsf/courses/dsr"
```

Are you happy with that place? if not, you should set your working directory to where all your data and script files are (or will be). Within RStudio you can go to 'Session' > 'Set working directory' > 'Choose directory'. Please do this now.

#### 4. Interactive introduction with swirl

Instead of clicking, you can use the function `setwd("/YOURPATH")`. For example, `setwd("/Users/MYNAME/MYFOLDER")` or `setwd("C:/Users/jenny/myrstuff")`. Make sure that the slashes are forward slashes and that you do not forget the apostrophes. R is case sensitive, so make sure you write capitals where necessary.

Whenever you want R to do something you need to use a function. It is like a command. All functions of R are organized in so-called packages or libraries. With the standard installation many packages are already installed. However, many more exist and some of them are really cool. For example, with `installed.packages()` all installed packages are listed. Or, with `swirl()`, you started swirl.

Of course, you can also go to the Packages window at the bottom right. If the box in front of the package name is ticked, the package is loaded (activated) and can be used. To see via Console which packages are loaded type in the console (`.packages()`)

```
(.packages())
```

```
[1] "stats"      "graphics"   "grDevices"  "utils"       "datasets"   "methods"  
[7] "base"
```

There are many more packages available on the R website. If you want to install and use a package (for example, the package called `geometry`) you should first install the package. Type `install.packages("geometry")` in the console. Don't be afraid about the many messages. Depending on your PC and your internet connection this may take some time.

```
install.packages("geometry")
```

After having installed a package, you need to load the package. That is a bit annoying but essential. Type in `library("geometry")` in the Console. You also did this for the swirl package (otherwise you couldn't have been doing these exercises).

```
library("geometry")
```

Check if the package is loaded typing `(.packages())`

```
(.packages())
```

Now, let's get started with the real programming.

R can be used as a calculator. You can just type your equation in the command window after the `>`. Type `10^2 + 36`.

```
10^2 + 36
```

```
[1] 136
```

And R gave the answer directly. By the way, spaces do not matter.

If you use brackets and forget to add the closing bracket, the `>` on the command line changes into a `+`. The `+` can also mean that R is still busy with some heavy computation. If you want R to quit what it was doing and give back the `>`, press ESC.

You can also give numbers a name. By doing so, they become so-called variables which can be used later. For example, you can type in the command window `A <- 4`.

```
A <- 4
```

#### 4. Interactive introduction with swirl

The `<-` is the so-called assignment operator. It allows you to assign data to a named object in order to store the data.

Don't be confused about the term object. All sorts of data are stored in so-called objects in R. All objects of a session are shown in the Environment window. In the second part of this course, I will introduce different data types.

You can see that `A` appeared in the environment window in the top right corner, which means that R now remembers what `A` is.

You can also ask R what `A` is. Just type `A` in the command window.

```
A
```

```
[1] 4
```

You can also do calculations with `A`. Type `A * 5`.

```
A * 5
```

```
[1] 20
```

If you specify `A` again, it will forget what value it had before. You can also assign a new value to `A` using the old one. Type `A <- A + 10`.

```
A <- A + 10
```

You can see that the value in the environment window changed.

To remove all variables from R's memory, type `rm(list=ls())`.

```
rm(list = ls())
```

You see that the environment window is now empty. You can also click the broom icon (`clear all`) in the environment window. You can see that RStudio then empties the environment window. If you only want to remove the variable `A`, you can type `rm(A)`.

Like in many other programs, R organizes numbers in scalars (a single number, 0-dimensional), vectors (a row of numbers, also called arrays, 1-dimensional) and matrices (like a table, 2-dimensional).

The `A` you defined before was a scalar. To define a vector with the numbers 3, 4 and 5, you need the function `c()`, which is short for concatenate (paste together). Type `B=c(3,4,5)`.

```
B <- c(3, 4, 5)
```

If you would like to compute the mean of all the elements in the vector `B` from the example above, you could type `(3+4+5)/3`. Try this

```
(3 + 4 + 5) / 3
```

```
[1] 4
```

But when the vector is very long, this is very boring and time-consuming work. This is why things you do often are automated in so-called functions. For example, type `mean(x=B)` and guess what this function `mean()` can do for you.

```
mean(x = B)
```

#### 4. Interactive introduction with swirl

```
[1] 4
```

Within the brackets you specify the arguments. Arguments give extra information to the function. In this case, the argument `x` says of which set of numbers (vector) the mean should be computed (namely of `B`). Sometimes, the name of the argument is not necessary; `mean(B)` works as well. Try it.

```
mean(B)
```

```
[1] 4
```

Compute the sum of 4, 5, 8 and 11 by first combining them into a vector and then using the function `sum`. Use the function `c` inside the function `sum`.

```
sum(c(4, 5, 8, 11))
```

```
[1] 28
```

The function `rnorm`, as another example, is a standard R function which creates random samples from a normal distribution. Type `rnorm(10)` and you will see 10 random numbers

```
rnorm(10)
```

```
[1] -0.60755509 -0.31113581 -0.12568079  0.76834476 -2.32547955 -0.18010883  
[7]  0.03819626  1.00268199  0.40764434  1.10742007
```

Here `rnorm` is the function and the 10 is an argument specifying how many random numbers you want - in this case 10 numbers (typing `n=10` instead of just 10 would also work). The result is 10 random numbers organised in a vector with length 10.

If you want 10 random numbers out of normal distribution with mean 1.2 and standard deviation 3.4 you can type `rnorm(10, mean=1.2, sd=3.4)`. Try this.

```
rnorm(10, mean = 1.2, sd = 3.4)
```

```
[1] -0.06122197  8.66105525 -2.81629111  1.07714911  3.59156880  3.31100577  
[7]  5.20839169  3.04920305 -1.70094353 -2.83630068
```

This shows that the same function (`rnorm()`) may have different interfaces and that R has so called named arguments (in this case `mean` and `sd`).

Comparing this example to the previous one also shows that for the function `rnorm` only the first argument (the number 10) is compulsory, and that R gives default values to the other so-called optional arguments. Use the help function to see which values are used as default by typing `?rnorm`.

```
?rnorm
```

You see the help page for this function in the help window on the right. RStudio has a nice features such as autocompletion and snapshots of the R documentation. For example, when you type `rnorm(` in the command window and press TAB, RStudio will show the possible arguments.

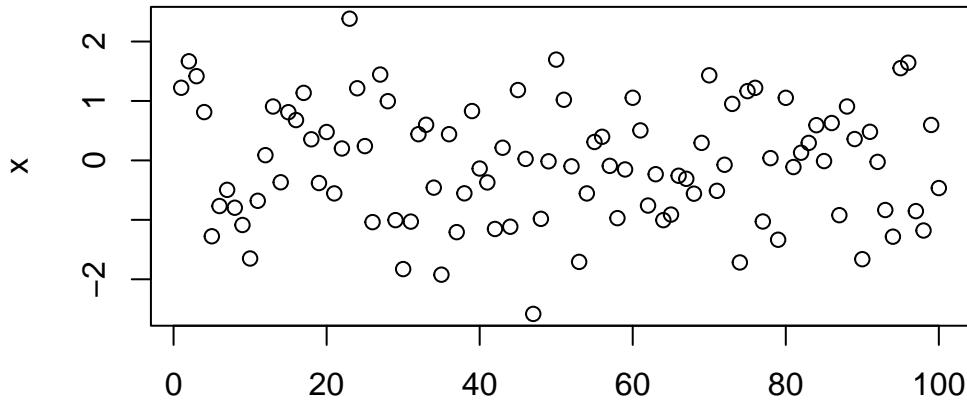
You can also store the output of the function in a variable. Type `x=rnorm(100)`.

```
x <- rnorm(100)
```

Now 100 random numbers are assigned to the variable x, which becomes a vector by this operation. You can see it appears in the Environment window.

R can also make graphs. Type `plot(x)` for a very simple example.

```
plot(x)
```



### Index

The 100 random numbers are now plotted in the plots window on the right.

You now are more familiar to RStudio and you know some basic R stuff. In particular, you know...

- ...that everything in R is said with functions,
- ...that functions can but don't have to have arguments,
- ...that you can install packages which contain functions,
- ...that you must load the installed packages every time you start a session in RStudio, and
- ...that this is just the beginning. Thus, please continue with the second module of this introduction.

After you have successfully finished learning module *huber-intro-1* please go ahead with the learning module *huber-intro-2* that is also part of my swirl course *swirl-it*.

## 4.3. swirl-it: huber-intro-2

**i** Click to see the full content of the module

Welcome to the second module. Again, if you find any errors or if you have suggestions for improvement, please let me know via [stephan.huber@hs-fresenius.de](mailto:stephan.huber@hs-fresenius.de) .

Before you start working, you should set your working directory to where all your data and script files are or should be stored. Within RStudio you can go to ‘Session’>‘Set working directory’, or you can type in `setwd(YOURPATH)`. Please do this now.

```
setwd("/home/sthu/Documents/mydir")
```

*Hint: Instead of clicking, you can also type `setwd("path")`, where you replace “path” with the location of your folder, for example `setwd("D:/R/swirl")`.*

R is an interpreter that uses a command line based environment. This means that you

#### 4. Interactive introduction with swirl

have to type commands, rather than use the mouse and menus. This has many advantages. Foremost, it is easy to get a full transcript of everything you did and you can replicate your work easy.

As already mentioned, all commands in R are functions where arguments come (or do not come) in round brackets after the function name.

You can store your workflow in files, the so-called scripts. These scripts have typically file names with the extension, e.g., foo.R .

You can open an editor window to edit these files by clicking ‘File’ and ‘New’. Try this. Under ‘File’ you also find the options ‘Open file...’, ‘Save’ and ‘Save as’. Alternatively, just type CTRL+SHIFT+N.

You can run (send to the Console window) part of the code by selecting lines and pressing CTRL+ENTER or click ‘Run’ in the editor window. If you do not select anything, R will run the line your cursor is on.

You can always run the whole script with the console command source, so e.g. for the script in the file foo.R you type source('foo.R'). You can also click ‘Run all’ in the editor window or type CTRL+SHIFT+S to run the whole script at once.

Make a script called firstscript.R. Therefore, open the editor window with ‘File’ > ‘New’. Type `plot(rnorm(100))` in the script, save it as `firstscript.R` in the working directory. Then type `source("firstscript.R")` on the command line.

```
source("firstscript.R")
```

Run your script again with `source("firstscript.R")`. The plot will change because new numbers are generated.

```
source("firstscript.R")
```

*Hint: Type `source("firstscript.R")` again or type `skip()` if you are not interested.*

Vectors were already introduced, but they can do more. Make a vector with numbers 1, 4, 6, 8, 10 and call it `vec1`.

*Hint: Type `vec1 <- c(1,4,6,8,10)`.*

```
vec1 <- c(1, 4, 6, 8, 10)
```

Elements in vectors can be addressed by standard [i] indexing. Select the 5th element of this vector by typing `vec1[5]`.

```
vec1[5]
```

Replace the 3rd element with a new number by typing `vec1[3]=12`.

```
vec1[3] <- 12
```

Ask R what the new version is of `vec1`.

```
vec1
```

You can also see the numbers of `vec1` in the environment window. Make a new vector `vec2` using the `seq()` (sequence) function by typing `seq(from=0, to=1, by=0.25)` and check its values in the environment window.

*Hint: Type `vec2 <- seq(from=0, to=1, by=0.25)`.*

#### 4. Interactive introduction with swirl

```
vec2 <- seq(from = 0, to = 1, by = 0.25)
```

Type `sum(vec1)`.

```
sum(vec1)
```

The function `sum` sums up the elements within a vector, leading to one number (a scalar). Now use `+` to add the two vectors.

*Hint: Type `vec1 + vec2`.*

```
vec1 + vec2
```

If you add two vectors of the same length, the first elements of both vectors are summed, and the second elements, etc., leading to a new vector of length 5 (just like in regular vector calculus).

Matrices are nothing more than 2-dimensional vectors. To define a matrix, use the function `matrix`. Make a matrix with `matrix(data=c(9,2,3,4,5,6),ncol=3)` and call it `mat`.

*Hint: Type `mat <- matrix(data=c(9,2,3,4,5,6),ncol=3)` or type `skip()` if you are not interested.*

```
mat <- matrix(data = c(9, 2, 3, 4, 5, 6), ncol = 3)
```

The third type of data structure treated here is the data frame. Time series are often ordered in data frames. A data frame is a matrix with names above the columns. This is nice, because you can call and use one of the columns without knowing in which position it is. Make a data frame with `t = data.frame(x = c(11,12,14), y = c(19,20,21), z = c(10,9,7))`.

```
t <- data.frame(x = c(11, 12, 14), y = c(19, 20, 21), z = c(10, 9, 7))
```

Ask R what `t` is.

*Hint: Type `t` or `skip()` if you are not interested.*

```
t
```

The data frame is called `t` and the columns have the names `x`, `y` and `z`. You can select one column by typing `t$z`. Try this.

```
t$z
```

Another option is to type `t[["z"]]`. Try this as well.

```
t[["z"]]
```

Compute the mean of column `z` in data frame `t`.

*Hint: Use function `mean` or type `skip()` if you are not interested.*

```
mean(t$z)
```

In the following question you will be asked to modify a script that will appear as soon as you move on from this question. When you have finished modifying the script, save your

#### 4. Interactive introduction with swirl

changes to the script and type `submit()` and the script will be evaluated. There will be some comments in the script that opens up. Be sure to read them!

Make a script file which constructs three random normal vectors of length 100. Call these vectors `x1`, `x2` and `x3`. Make a data frame called `t` with three columns (called `a`, `b` and `c`) containing respectively `x1`, `x1+x2` and `x1+x2+x3`. Call `plot(t)` for this data frame. Then, save it and type `submit()` on the command line.

*Hint: Type `plot(rnorm(100))` in the script, save it and type `submit()` on the command line.*

```
# Text behind the #-sign is not evaluated as code by R.  
# This is useful, because it allows you to add comments explaining what the script does.  
  
# In this script, replace the ... with the appropriate commands.  
  
x1 <- ...  
x2 <- ...  
x3 <- ...  
t <- ...  
plot(...)
```

#### Result

```
# Text behind the #-sign is not evaluated as code by R.  
# This is useful, because it allows you to add comments explaining what the script does.  
  
# In this script, replace the ... with the appropriate commands.  
  
x1 <- rnorm(100)  
x2 <- rnorm(100)  
x3 <- rnorm(100)  
t <- data.frame(a = x1, b = x1 + x2, c = x1 + x2 + x3)  
plot(t)
```

Do you understand the results?

Another basic structure in R is a list. The main advantage of lists is that the `columns` (they are not really ordered in columns any more, but are more a collection of vectors) don't have to be of the same length, unlike matrices and data frames. Make this list `L <- list(one=1, two=c(1,2), five=seq(0, 1, length=5))`.

```
L <- list(one = 1, two = c(1, 2), five = seq(0, 1, length = 5))
```

The list `L` has names and values. You can type `L` to see the contents.

```
L
```

`L` also appeared in the environment window. To find out what's in the list, type `names(L)`.

```
names(L)
```

Add 10 to the column called `five`.

*Hint: Type `L$five + 10`*

#### 4. Interactive introduction with swirl

```
L$five + 10
```

Plotting is an important statistical activity. So it should not come as a surprise that R has many plotting facilities. Type `plot(rnorm(100), type="l", col="gold")`.

*Hint: The symbol between quotes after the `type=`, is the letter `l`, not the number 1. To see the result you can also just type `skip()`.*

```
plot(rnorm(100), type = "l", col = "gold")
```

Hundred random numbers are plotted by connecting the points by lines in a gold color. Another very simple example is the classical statistical histogram plot, generated by the simple command `hist`. Make a histogram of 100 random numbers.

*Hint: Type `hist(rnorm(100))`*

```
hist(rnorm(100))
```

The script that opens up is the same as the script you made before, but with more plotting commands. Type `submit()` on the command line to run it (you don't have to change anything yet).

*Hint: Change plotting parameters in the script, save it and type `submit()` on the command line.*

```
# Text behind the #-sign is not evaluated as code by R.  
# This is useful, because it allows you to add comments explaining what the script does.  
  
# Make data frame  
x1 <- rnorm(100)  
x2 <- rnorm(100)  
x3 <- rnorm(100)  
t <- data.frame(a = x1, b = x1 + x2, c = x1 + x2 + x3)  
  
# Plot data frame  
plot(t$a, type = "l", ylim = range(t), lwd = 3, col = rgb(1, 0, 0, 0.3))  
lines(t$b, type = "s", lwd = 2, col = rgb(0.3, 0.4, 0.3, 0.9))  
points(t$c, pch = 20, cex = 4, col = rgb(0, 0, 1, 0.3))  
  
# Note that with plot you get a new plot window while points and lines add to the previous
```

Try to find out by experimenting what the meaning is of `rgb`, the last argument of `rgb`, `lwd`, `pch`, `cex`. Type `play()` on the command line to experiment. Modify lines 11, 12 and 13 of the script by putting your cursor there and pressing CTRL+ENTER. When you are finished, type `nxt()` and then `?par`.

*Hint: Type `?par` or type `skip()` if you are not interested.*

```
?par
```

You searched for `par` in the R help. This is a useful page to learn more about formatting plots. Google 'R color chart' for a pdf file with a wealth of color options.

To copy your plot to a document, go to the plots window, click the 'Export' button, choose the nicest width and height and click 'Copy' or 'Save'.

After having almost completed the second learning module, you are getting closer to become a nerd as you know...

...that everything in R is stored in objects (values, vectors, matrices, lists, or data frames),  
...that you should always work in scripts and send code from scripts to the Console,  
...that you can do it if you don't give up.

Please continue choosing another `swirl` learning module.

## 4.4. `swirl-it`: Data analytical basics

In my `swirl` modules *huber-data-1*, *huber-data-2*, and *huber-data-3* I introduce some very basic statistical principles on how to analyse data.

## 4.5. `swirl-it`: The `tidyverse` package

I compiled a short `swirl` module to introduce the *tidyverse* universe. This is a powerful collection of packages which I discuss later on. The learning module is also part of my *swirl-it* course.

## 4.6. Other `swirl` modules

You can also install some other courses. You find a list of courses here <http://swirlstats.com/scn/index.html> or here [https://github.com/swirldev/swirl\\_courses](https://github.com/swirldev/swirl_courses).

I recommend this one as it gives a general overview on very basic principles of R:

```
library(swirl)
install_course_github("swirldev", "R_Programming_E")
swirl()
```

# 5. Kickstart

Ever got a kick that actually moved you forward? Well, let's kickstart your R adventure by walking you through a typical data analysis workflow in R, covering everything from setting up your environment to performing data analysis and visualization. Along the way, we'll also tackle some common troubleshooting to smooth out any bumps in the road. Please don't worry if you don't understand some lines of code. You will learn that later on, however, you hopefully will get a sense of what it is like to work with a command line based program.

## 5.1. Analysing the association of weight and the price of cars

Before we start, we need to ensure that all necessary libraries are installed and loaded. We use the pacman package for convenient package management.

```
# Install and load required libraries
# Installs 'pacman' if not already available, which is used for package management
if (!require(pacman)) install.packages("pacman")

# Unload all previously loaded packages to start fresh
suppressMessages(pacman::p_unload(all))

# Load necessary packages for data manipulation, cleaning, and visualization
pacman::p_load(
  tidyverse, # A suite of packages designed for data science that includes tools for data ma
  haven,    # Used for importing and exporting data with SPSS, Stata, and SAS formats.
  janitor,  # Provides functions for examining and cleaning data, such as `clean_names()` `a
  WDI,      # Facilitates downloading data from the World Bank's World Development Indicate
  wbstats   # Provides an interface to the World Bank's APIs for a comprehensive range of d
)

# Set the working directory to a project-specific folder
setwd("~/Dropbox/hsf/courses/dsr")

# Clear the current environment of any objects
rm(list = ls())
```

Now, let us load the dataset from a Stata file (`auto.dta`) and explore its basic properties.

```
# Load data from a Stata file available online
auto <- read_dta("http://www.stata-press.com/data/r18/auto.dta")
# 'auto': Dataset contains information about different car models

# Display basic information about the dataset
ncol(auto) # Number of columns
```

## 5. Kickstart

```
[1] 12
```

```
nrow(auto) # Number of rows
```

```
[1] 74
```

```
dim(auto) # Dimensions of the dataset
```

```
[1] 74 12
```

```
names(auto) # Names of variables
```

```
[1] "make"          "price"         "mpg"           "rep78"          "headroom"  
[6] "trunk"         "weight"        "length"        "turn"           "displacement"  
[11] "gear_ratio"    "foreign"
```

```
head(auto) # First few rows
```

```
# A tibble: 6 x 12  
  make      price   mpg rep78 headroom trunk weight length turn displacement  
  <chr>     <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>  
1 AMC Concord  4099    22     3    2.5    11   2930    186    40      121  
2 AMC Pacer    4749    17     3     3    11   3350    173    40      258  
3 AMC Spirit   3799    22    NA     3    12   2640    168    35      121  
4 Buick Centur~ 4816    20     3    4.5    16   3250    196    40      196  
5 Buick Electr~ 7827    15     4     4    20   4080    222    43      350  
6 Buick LeSab~  5788    18     3     4    21   3670    218    43      231  
# i 2 more variables: gear_ratio <dbl>, foreign <dbl+lbl>
```

```
tail(auto) # Last few rows
```

```
# A tibble: 6 x 12  
  make      price   mpg rep78 headroom trunk weight length turn displacement  
  <chr>     <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>  
1 Toyota Coro~  5719    18     5     2    11   2670    175    36      134  
2 VW Dasher    7140    23     4    2.5    12   2160    172    36      97  
3 VW Diesel    5397    41     5     3    15   2040    155    35      90  
4 VW Rabbit    4697    25     4     3    15   1930    155    35      89  
5 VW Scirocco  6850    25     4     2    16   1990    156    36      97  
6 Volvo 260    11995   17     5    2.5    14   3170    193    37      163  
# i 2 more variables: gear_ratio <dbl>, foreign <dbl+lbl>
```

```
summary(auto) # Summary statistics for each column
```

## 5. Kickstart

```

      make          price         mpg        rep78
Length:74      Min.   : 3291   Min.   :12.00   Min.   :1.000
Class :character  1st Qu.: 4220   1st Qu.:18.00   1st Qu.:3.000
Mode  :character  Median : 5006   Median :20.00   Median :3.000
                  Mean    : 6165   Mean    :21.30   Mean    :3.406
                  3rd Qu.: 6332   3rd Qu.:24.75   3rd Qu.:4.000
                  Max.    :15906   Max.    :41.00   Max.    :5.000
                               NA's    :5

      headroom       trunk        weight       length       turn
Min.   :1.500   Min.   : 5.00   Min.   :1760   Min.   :142.0   Min.   :31.00
1st Qu.:2.500  1st Qu.:10.25  1st Qu.:2250  1st Qu.:170.0  1st Qu.:36.00
Median :3.000  Median :14.00  Median :3190  Median :192.5  Median :40.00
Mean   :2.993  Mean   :13.76  Mean   :3019  Mean   :187.9  Mean   :39.65
3rd Qu.:3.500  3rd Qu.:16.75  3rd Qu.:3600  3rd Qu.:203.8  3rd Qu.:43.00
Max.   :5.000  Max.   :23.00  Max.   :4840  Max.   :233.0  Max.   :51.00

      displacement   gear_ratio   foreign
Min.   : 79.0   Min.   :2.190   Min.   :0.0000
1st Qu.:119.0  1st Qu.:2.730  1st Qu.:0.0000
Median :196.0   Median :2.955  Median :0.0000
Mean   :197.3   Mean   :3.015  Mean   :0.2973
3rd Qu.:245.2  3rd Qu.:3.352  3rd Qu.:1.0000
Max.   :425.0   Max.   :3.890  Max.   :1.0000

```

```
glimpse(auto) # Compact display of the structure of the dataset
```

```

Rows: 74
Columns: 12
$ make           <chr> "AMC Concord", "AMC Pacer", "AMC Spirit", "Buick Century"~
$ price          <dbl> 4099, 4749, 3799, 4816, 7827, 5788, 4453, 5189, 10372, 40~
$ mpg            <dbl> 22, 17, 22, 20, 15, 18, 26, 20, 16, 19, 14, 14, 21, 29, 1~
$ rep78          <dbl> 3, 3, NA, 3, 4, 3, NA, 3, 3, 3, 2, 3, 3, 4, 3, 2, 2, 3~
$ headroom        <dbl> 2.5, 3.0, 3.0, 4.5, 4.0, 4.0, 3.0, 2.0, 3.5, 3.5, 4.0, 3.~
$ trunk           <dbl> 11, 11, 12, 16, 20, 21, 10, 16, 17, 13, 20, 16, 13, 9, 20~
$ weight          <dbl> 2930, 3350, 2640, 3250, 4080, 3670, 2230, 3280, 3880, 340~
$ length          <dbl> 186, 173, 168, 196, 222, 218, 170, 200, 207, 200, 221, 20~
$ turn             <dbl> 40, 40, 35, 40, 43, 43, 34, 42, 43, 42, 44, 43, 45, 34, 4~
$ displacement    <dbl> 121, 258, 121, 196, 350, 231, 304, 196, 231, 231, 425, 35~
$ gear_ratio      <dbl> 3.58, 2.53, 3.08, 2.93, 2.41, 2.73, 2.87, 2.93, 2.93, 3.0~
$ foreign         <dbl+lbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~

```

```
print(auto, n = Inf) # Print all rows of the dataset
```

```

# A tibble: 74 x 12
  make      price     mpg rep78 headroom trunk weight length turn displacement
  <chr>     <dbl>   <dbl> <dbl>    <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
1 AMC Concord 4099     22     3     2.5     11    2930     186     40     121
2 AMC Pacer   4749     17     3     3       11    3350     173     40     258
3 AMC Spirit   3799     22    NA     3       12    2640     168     35     121

```

## 5. Kickstart

4 Buick Cent~	4816	20	3	4.5	16	3250	196	40	196
5 Buick Elec~	7827	15	4	4	20	4080	222	43	350
6 Buick LeSa~	5788	18	3	4	21	3670	218	43	231
7 Buick Opel	4453	26	NA	3	10	2230	170	34	304
8 Buick Regal	5189	20	3	2	16	3280	200	42	196
9 Buick Rivi~	10372	16	3	3.5	17	3880	207	43	231
10 Buick Skyl~	4082	19	3	3.5	13	3400	200	42	231
11 Cad. Devil~	11385	14	3	4	20	4330	221	44	425
12 Cad. Eldor~	14500	14	2	3.5	16	3900	204	43	350
13 Cad. Sevil~	15906	21	3	3	13	4290	204	45	350
14 Chev. Chev~	3299	29	3	2.5	9	2110	163	34	231
15 Chev. Impa~	5705	16	4	4	20	3690	212	43	250
16 Chev. Mali~	4504	22	3	3.5	17	3180	193	31	200
17 Chev. Mont~	5104	22	2	2	16	3220	200	41	200
18 Chev. Monza	3667	24	2	2	7	2750	179	40	151
19 Chev. Nova	3955	19	3	3.5	13	3430	197	43	250
20 Dodge Colt	3984	30	5	2	8	2120	163	35	98
21 Dodge Dipl~	4010	18	2	4	17	3600	206	46	318
22 Dodge Magn~	5886	16	2	4	17	3600	206	46	318
23 Dodge St. ~	6342	17	2	4.5	21	3740	220	46	225
24 Ford Fiesta	4389	28	4	1.5	9	1800	147	33	98
25 Ford Musta~	4187	21	3	2	10	2650	179	43	140
26 Linc. Cont~	11497	12	3	3.5	22	4840	233	51	400
27 Linc. Mark~	13594	12	3	2.5	18	4720	230	48	400
28 Linc. Vers~	13466	14	3	3.5	15	3830	201	41	302
29 Merc. Bobc~	3829	22	4	3	9	2580	169	39	140
30 Merc. Coug~	5379	14	4	3.5	16	4060	221	48	302
31 Merc. Marq~	6165	15	3	3.5	23	3720	212	44	302
32 Merc. Mona~	4516	18	3	3	15	3370	198	41	250
33 Merc. XR-7	6303	14	4	3	16	4130	217	45	302
34 Merc. Zeph~	3291	20	3	3.5	17	2830	195	43	140
35 Olds 98	8814	21	4	4	20	4060	220	43	350
36 Olds Cutf ~	5172	19	3	2	16	3310	198	42	231
37 Olds Cutla~	4733	19	3	4.5	16	3300	198	42	231
38 Olds Delta~	4890	18	4	4	20	3690	218	42	231
39 Olds Omega	4181	19	3	4.5	14	3370	200	43	231
40 Olds Starf~	4195	24	1	2	10	2730	180	40	151
41 Olds Toron~	10371	16	3	3.5	17	4030	206	43	350
42 Plym. Arrow	4647	28	3	2	11	3260	170	37	156
43 Plym. Champ	4425	34	5	2.5	11	1800	157	37	86
44 Plym. Hori~	4482	25	3	4	17	2200	165	36	105
45 Plym. Sapp~	6486	26	NA	1.5	8	2520	182	38	119
46 Plym. Vola~	4060	18	2	5	16	3330	201	44	225
47 Pont. Cata~	5798	18	4	4	20	3700	214	42	231
48 Pont. Fire~	4934	18	1	1.5	7	3470	198	42	231
49 Pont. Gran~	5222	19	3	2	16	3210	201	45	231
50 Pont. Le M~	4723	19	3	3.5	17	3200	199	40	231
51 Pont. Phoe~	4424	19	NA	3.5	13	3420	203	43	231
52 Pont. Sunb~	4172	24	2	2	7	2690	179	41	151
53 Audi 5000	9690	17	5	3	15	2830	189	37	131
54 Audi Fox	6295	23	3	2.5	11	2070	174	36	97

## 5. Kickstart

```

55 BMW 320i    9735    25     4     2.5    12   2650    177    34    121
56 Datsun 200   6229    23     4     1.5     6   2370    170    35    119
57 Datsun 210   4589    35     5     2      8   2020    165    32     85
58 Datsun 510   5079    24     4     2.5     8   2280    170    34    119
59 Datsun 810   8129    21     4     2.5     8   2750    184    38    146
60 Fiat Strada  4296    21     3     2.5    16   2130    161    36    105
61 Honda Acco~  5799    25     5     3      10   2240    172    36    107
62 Honda Civic   4499    28     4     2.5     5   1760    149    34     91
63 Mazda GLC    3995    30     4     3.5    11   1980    154    33     86
64 Peugeot 604   12990   14     NA    3.5    14   3420    192    38    163
65 Renault Le~   3895    26     3     3      10   1830    142    34     79
66 Subaru        3798    35     5     2.5    11   2050    164    36     97
67 Toyota Cel~   5899    18     5     2.5    14   2410    174    36    134
68 Toyota Cor~   3748    31     5     3      9    2200    165    35     97
69 Toyota Cor~   5719    18     5     2      11   2670    175    36    134
70 VW Dasher    7140    23     4     2.5    12   2160    172    36     97
71 VW Diesel    5397    41     5     3      15   2040    155    35     90
72 VW Rabbit    4697    25     4     3      15   1930    155    35     89
73 VW Scirocco  6850    25     4     2      16   1990    156    36     97
74 Volvo 260    11995   17     5     2.5    14   3170    193    37    163
# i 2 more variables: gear_ratio <dbl>, foreign <dbl+lbl>

```

The data seems to be a cross-section of cars. Let us check if the variable `make` identifies each line uniquely:

```

# Check for duplicate entries based on the 'make' variable
auto |>
  get_dupes(make)

```

No duplicate combinations found of: make

```

# A tibble: 0 x 13
# i 13 variables: make <chr>, dupe_count <int>, price <dbl>, mpg <dbl>,
#   rep78 <dbl>, headroom <dbl>, trunk <dbl>, weight <dbl>, length <dbl>,
#   turn <dbl>, displacement <dbl>, gear_ratio <dbl>, foreign <dbl+lbl>

```

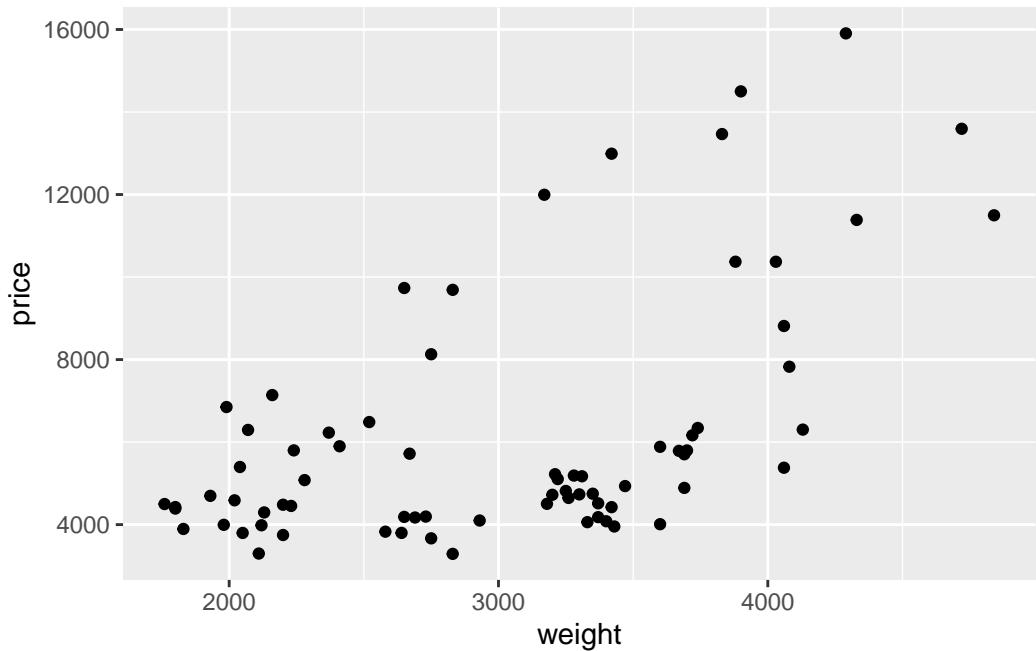
Indeed, the variable `make` has no duplicates. Now, let's make and save some graphical visualizations:

```

# Create and display a scatter plot of car price versus weight
plot_weight_price <- ggplot(auto, aes(x = weight, y = price)) +
  geom_point()
plot_weight_price

```

## 5. Kickstart

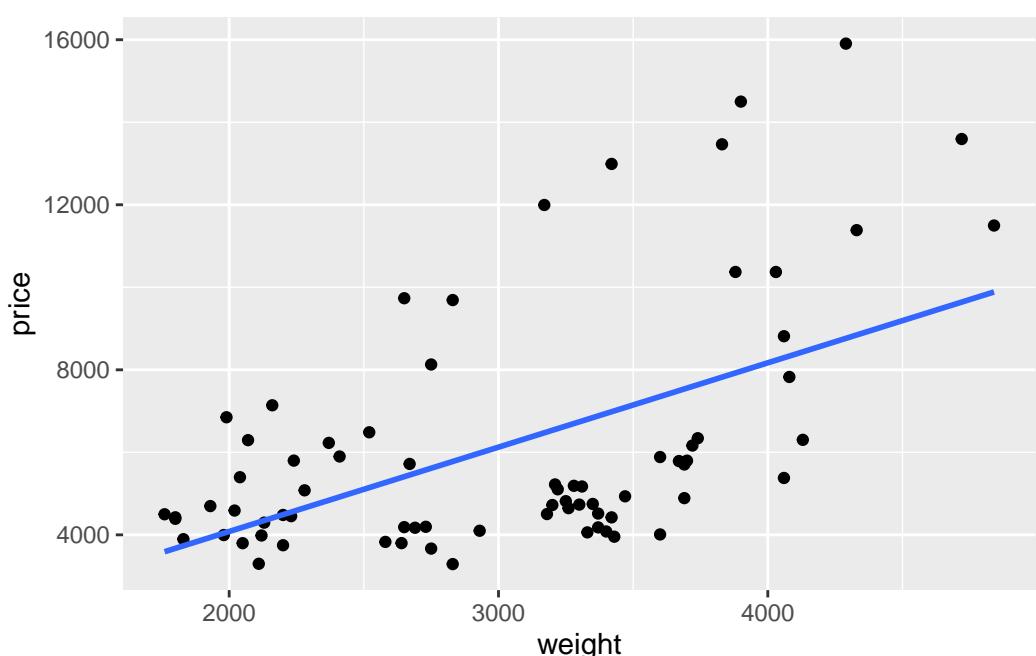


```
# Save the plot to a file
ggsave("fig/plot_weight_price.png", plot = plot_weight_price, dpi = 300)
```

Saving 5.5 x 3.5 in image

```
# Create a scatter plot with a linear regression line of price vs weight
plot_weight_price_fit <- ggplot(auto, aes(x = weight, y = price)) +
  geom_point() +
  geom_smooth(method = lm, se = FALSE) # 'lm' denotes linear model, 'se' is standard error
plot_weight_price_fit
```

`geom\_smooth()` using formula = 'y ~ x'



## 5. Kickstart

```
# Save the plot to a file
ggsave("fig/plot_weight_price_fit.png", plot = plot_weight_price_fit, dpi = 300)
```

```
Saving 5.5 x 3.5 in image
`geom_smooth()` using formula = 'y ~ x'
```

Let us perform a linear regression to quantify the impact of `weight` on `price`:

```
# Perform a linear regression to analyze the relationship between weight and price
reg_result <- lm(price ~ weight , data = auto)
summary(reg_result) # Display the regression results
```

Call:

```
lm(formula = price ~ weight, data = auto)
```

Residuals:

Min	1Q	Median	3Q	Max
-3341.9	-1828.3	-624.1	1232.1	7143.7

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-6.7074	1174.4296	-0.006	0.995
weight	2.0441	0.3768	5.424	7.42e-07 ***
---				
Signif. codes:	0	'***'	0.001	'**'
			0.01	'*'
			0.05	'.'
			0.1	' '
			1	

Residual standard error: 2502 on 72 degrees of freedom

Multiple R-squared: 0.2901, Adjusted R-squared: 0.2802

F-statistic: 29.42 on 1 and 72 DF, p-value: 7.416e-07

## 5.2. Accessing World Bank's *World Development Indicators*

The World Wide Web is a treasure trove of data, and most major databases offer researchers direct download options. Numerous user-supplied packages are available for seamless access to such data. In this section, I present two popular packages that facilitate the downloading of data from the World Bank's *World Development Indicators*:

**WDI (World Development Indicators) Package** - Official CRAN package documentation: [WDI on CRAN](#) - Source code on GitHub: [WDI GitHub Repository](#)

**wbstats (World Bank Statistics) Package** - Official CRAN package documentation: [wbstats on CRAN](#) - Source code on GitHub: [wbstats GitHub Repository](#)

Now, let's download some GDP data and explore how to manipulate it. This exercise will demonstrate practical applications of the tools.

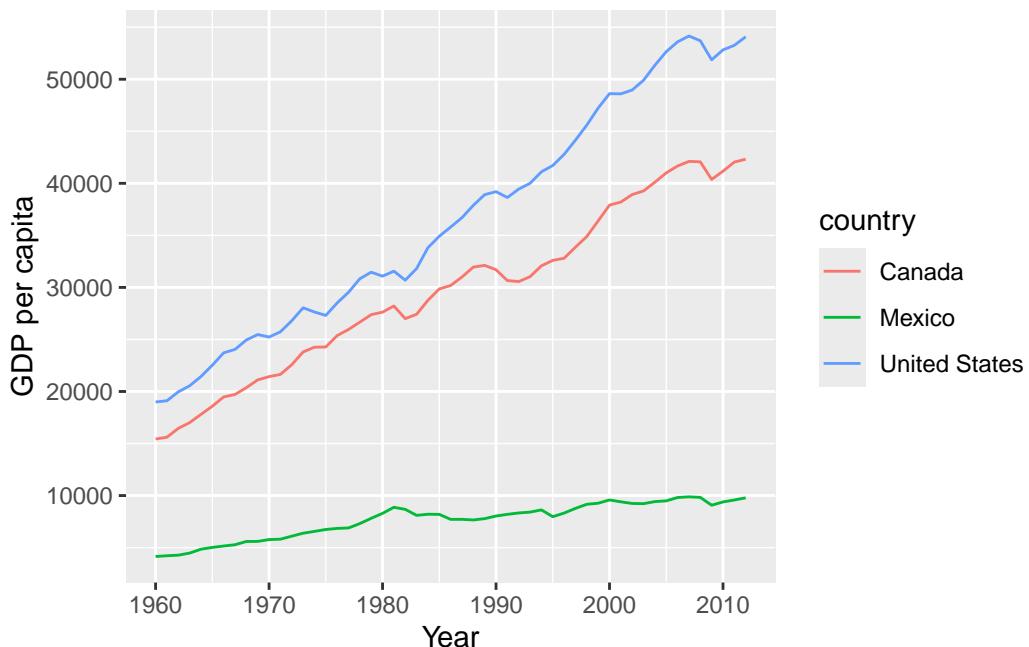
```
# Search for GDP indicators and display the first 10
WDIsearch("gdp") [1:10, ]
```

## 5. Kickstart

	indicator	name
689	6.0.GDP_current	GDP (current \$)
690	6.0.GDP_growth	GDP growth (annual %)
691	6.0.GDP_usd	GDP (constant 2005 \$)
692	6.0.GDPpc_constant	GDP per capita, PPP (constant 2011 international \$)
2096	BG.GSR.NFSV.GD.ZS	Trade in services (% of GDP)
2097	BG.KAC.FNEI.GD.PP.ZS	Gross private capital flows (% of GDP, PPP)
2098	BG.KAC.FNEI.GD.ZS	Gross private capital flows (% of GDP)
2099	BG.KLT.DINV.GD.PP.ZS	Gross foreign direct investment (% of GDP, PPP)
2100	BG.KLT.DINV.GD.ZS	Gross foreign direct investment (% of GDP)
2401	BI.WAG.TOTL.GD.ZS	Wage bill as a percentage of GDP

```
# Retrieve GDP per capita data for specified countries and years
df_WDI <- WDI(
  indicator = "NY.GDP.PCAP.KD",
  country = c("MX", "CA", "US"),
  start = 1960,
  end = 2012
)

# Plot GDP per capita over time for the specified countries
ggplot(df_WDI, aes(year, NY.GDP.PCAP.KD, color = country)) +
  geom_line() +
  xlab("Year") +
  ylab("GDP per capita")
```

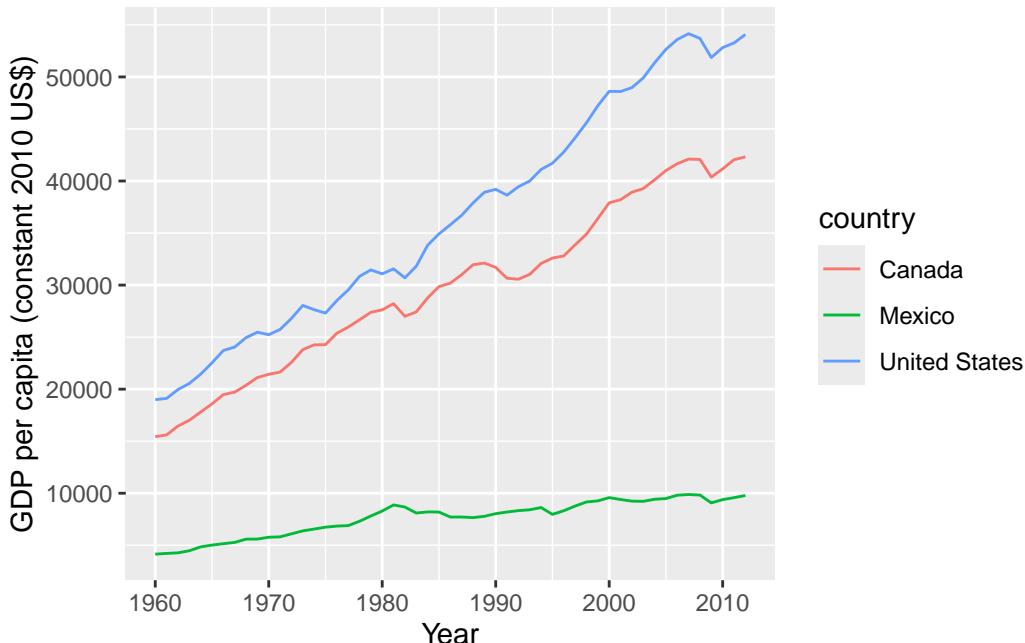


```
# Retrieve GDP per capita data for specified countries and years using the wbstats package
df_wb <- wb_data(
  indicator = "NY.GDP.PCAP.KD",
  country = c("MX", "CA", "US"),
  start = 1960,
  end = 2012,
```

## 5. Kickstart

```
    return_wide = TRUE
)

# Plot GDP per capita over time for the specified countries
ggplot(df_wb, aes(date, NY.GDP.PCAP.KD, color = country)) +
  geom_line() +
  xlab("Year") +
  ylab("GDP per capita (constant 2010 US$)")
```



The latter graph appears empty. Why? Let's take a closer look at the data to identify any discrepancies that might explain this issue:

```
# Look at the data types year and date are different:
glimpse(df_WDI)
```

```
Rows: 159
Columns: 5
$ country      <chr> "Canada", "Canada", "Canada", "Canada", "Canada", "Cana~
$ iso2c        <chr> "CA", "CA", "CA", "CA", "CA", "CA", "CA", "CA", "CA", "~
$ iso3c        <chr> "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", ~
$ year         <int> 2012, 2011, 2010, 2009, 2008, 2007, 2006, 2005, 2004, 2~
$ NY.GDP.PCAP.KD <dbl> 42320.64, 42043.64, 41164.34, 40376.42, 42067.57, 42106~
```

```
glimpse(df_wb)
```

```
Rows: 159
Columns: 9
$ iso2c        <chr> "CA", "CA", "CA", "CA", "CA", "CA", "CA", "CA", "CA", "~
$ iso3c        <chr> "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", ~
$ country      <chr> "Canada", "Canada", "Canada", "Canada", "Canada", "Cana~
$ date         <dbl> 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968, 1~
```

## 5. Kickstart

```
$ NY.GDP.PCAP.KD <dbl> 15432.47, 15605.52, 16455.75, 17007.69, 17800.86, 18585~  
$ unit <chr> NA, NA,~  
$ obs_status <chr> NA, NA,~  
$ footnote <chr> NA, NA,~  
$ last_updated <date> 2025-04-15, 2025-04-15, 2025-04-15, 2025-04-15, 2025-0~
```

The answer is: A lineplot with a character variable (date is <chr>) on the x-axis does not work!

Now, let us manipulate the df\_wb data so that the two dataset are equal:

```
df_wb_cln <- df_wb |>  
  # Convert 'date' in df_wb from character to integer  
  mutate(year = as.integer(date)) |>  
  # Since 'year' has been created, remove the original 'date' column  
  select(-date) |>  
  # Relocate columns to organize the data frame  
  relocate(country, iso2c, iso3c, year, NY.GDP.PCAP.KD)  
  
glimpse(df_WDI)
```

```
Rows: 159  
Columns: 5  
$ country <chr> "Canada", "Canada", "Canada", "Canada", "Canada", "Cana~  
$ iso2c <chr> "CA", ~  
$ iso3c <chr> "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", ~  
$ year <int> 2012, 2011, 2010, 2009, 2008, 2007, 2006, 2005, 2004, 2~  
$ NY.GDP.PCAP.KD <dbl> 42320.64, 42043.64, 41164.34, 40376.42, 42067.57, 42106~
```

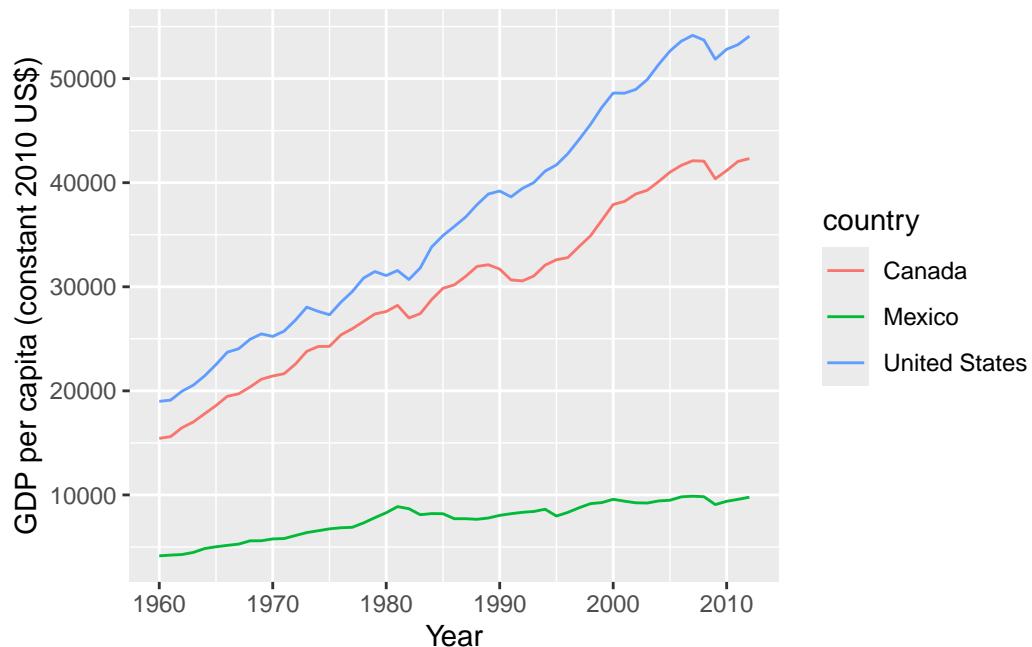
```
glimpse(df_wb_cln)
```

```
Rows: 159  
Columns: 9  
$ country <chr> "Canada", "Canada", "Canada", "Canada", "Canada", "Cana~  
$ iso2c <chr> "CA", ~  
$ iso3c <chr> "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", ~  
$ year <int> 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968, 1~  
$ NY.GDP.PCAP.KD <dbl> 15432.47, 15605.52, 16455.75, 17007.69, 17800.86, 18585~  
$ unit <chr> NA, NA,~  
$ obs_status <chr> NA, NA,~  
$ footnote <chr> NA, NA,~  
$ last_updated <date> 2025-04-15, 2025-04-15, 2025-04-15, 2025-04-15, 2025-0~
```

Now it works:

```
# Plot GDP per capita over time for the specified countries  
ggplot(df_wb_cln, aes(year, NY.GDP.PCAP.KD, color = country)) +  
  geom_line() +  
  xlab("Year") +  
  ylab("GDP per capita (constant 2010 US$)")
```

## 5. Kickstart



### Solution

The script uses the following functions: `aes`, `as.integer`, `c`, `dim`, `geom_line`, `geom_point`, `geom_smooth`, `get_dupes`, `ggplot`, `ggsave`, `glimpse`, `head`, `lm`, `mutate`, `names`, `nrow`, `print`, `read_dta`, `relocate`, `select`, `setwd`, `summary`, `tail`, `wb`, `WDI`, `WDIsearch`, `xlab`, `ylab`.

### R script

## 5. Kickstart

```
# This script demonstrates a typical data analysis workflow in R
# ----

# Install and load required libraries
# Installs 'pacman' if not already available, which is used for package management
if (!require(pacman)) install.packages("pacman")

# Unload all previously loaded packages to start fresh
suppressMessages(pacman::p_unload(all))

# Load necessary packages for data manipulation, cleaning, and visualization
pacman::p_load(
    tidyverse, # A suite of packages designed for data science that includes tools for
    haven, # Used for importing and exporting data with SPSS, Stata, and SAS formats.
    janitor, # Provides functions for examining and cleaning data, such as `clean_names`
    WDI, # Facilitates downloading data from the World Bank's World Development Indicators
    wbstats # Provides an interface to the World Bank's APIs for a comprehensive range
)

# Set the working directory to a project-specific folder
setwd("~/Dropbox/hsf/courses/dsr")

# Clear the current environment of any objects
rm(list = ls())

# ----

# Load data from a Stata file available online
auto <- read_dta("http://www.stata-press.com/data/r8/auto.dta")
# 'auto': Dataset contains information about different car models

# Display basic information about the dataset
ncol(auto) # Number of columns
nrow(auto) # Number of rows
dim(auto) # Dimensions of the dataset
names(auto) # Names of variables
head(auto) # First few rows
tail(auto) # Last few rows
summary(auto) # Summary statistics for each column
glimpse(auto) # Compact display of the structure of the dataset
print(auto, n = Inf) # Print all rows of the dataset

# Check for duplicate entries based on the 'make' variable
auto |>
    get_dupes(make)

# Create and display a scatter plot of car price versus weight
plot_weight_price <- ggplot(auto, aes(x = weight, y = price)) +
    geom_point()
plot_weight_price

# Save the plot to a file
ggsave("fig/plot_weight_price.png", plot = plot_weight_price, dpi = 300)

# Create a scatter plot with a linear regression line of price vs weight
plot_weight_price_fit <- ggplot(auto, aes(x = weight, y = price)) +
```

## Output of the R script

```
# This script demonstrates a typical data analysis workflow in R
# ----

# Install and load required libraries
# Installs 'pacman' if not already available, which is used for package management
if (!require(pacman)) install.packages("pacman")

# Unload all previously loaded packages to start fresh
suppressMessages(pacman::p_unload(all))

# Load necessary packages for data manipulation, cleaning, and visualization
pacman::p_load(
  tidyverse, # A suite of packages designed for data science that includes tools for
  haven, # Used for importing and exporting data with SPSS, Stata, and SAS formats.
  janitor, # Provides functions for examining and cleaning data, such as `clean_names`
  WDI, # Facilitates downloading data from the World Bank's World Development Indicators
  wbstats # Provides an interface to the World Bank's APIs for a comprehensive range
)

# Set the working directory to a project-specific folder
setwd("~/Dropbox/hsf/courses/dsr")

# Clear the current environment of any objects
rm(list = ls())

# ----

# Load data from a Stata file available online
auto <- read_dta("http://www.stata-press.com/data/r8/auto.dta")
# 'auto': Dataset contains information about different car models

# Display basic information about the dataset
ncol(auto) # Number of columns
```

[1] 12

```
nrow(auto) # Number of rows
```

[1] 74

```
dim(auto) # Dimensions of the dataset
```

[1] 74 12

```
names(auto) # Names of variables
```

## 5. Kickstart

```
[1] "make"          "price"         "mpg"           "rep78"        "headroom"
[6] "trunk"         "weight"        "length"        "turn"         "displacement"
[11] "gear_ratio"   "foreign"

head(auto) # First few rows

# A tibble: 6 x 12
  make      price    mpg rep78 headroom trunk weight length turn displacement
  <chr>     <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 AMC Concord 4099    22     3    2.5    11  2930    186    40      121
2 AMC Pacer   4749    17     3     3    11  3350    173    40      258
3 AMC Spirit  3799    22    NA     3    12  2640    168    35      121
4 Buick Centur~ 4816    20     3    4.5    16  3250    196    40      196
5 Buick Elect~ 7827    15     4     4    20  4080    222    43      350
6 Buick LeSab~ 5788    18     3     4    21  3670    218    43      231
# i 2 more variables: gear_ratio <dbl>, foreign <dbl+lbl>

tail(auto) # Last few rows

# A tibble: 6 x 12
  make      price    mpg rep78 headroom trunk weight length turn displacement
  <chr>     <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 Toyota Coro~ 5719    18     5     2    11  2670    175    36      134
2 VW Dasher   7140    23     4    2.5    12  2160    172    36      97
3 VW Diesel   5397    41     5     3    15  2040    155    35      90
4 VW Rabbit   4697    25     4     3    15  1930    155    35      89
5 VW Scirocco 6850    25     4     2    16  1990    156    36      97
6 Volvo 260   11995   17     5    2.5    14  3170    193    37      163
# i 2 more variables: gear_ratio <dbl>, foreign <dbl+lbl>

summary(auto) # Summary statistics for each column

  make            price          mpg          rep78
Length:74      Min.   :3291   Min.   :12.00   Min.   :1.000
Class :character 1st Qu.:4220   1st Qu.:18.00   1st Qu.:3.000
Mode  :character Median :5006   Median :20.00   Median :3.000
                  Mean   :6165   Mean   :21.30   Mean   :3.406
                  3rd Qu.:6332   3rd Qu.:24.75   3rd Qu.:4.000
                  Max.  :15906   Max.  :41.00   Max.  :5.000
                  NA's   :5
  headroom        trunk          weight         length        turn
Min.   :1.500   Min.   :5.00   Min.   :1760   Min.   :142.0   Min.   :31.00
1st Qu.:2.500  1st Qu.:10.25  1st Qu.:2250  1st Qu.:170.0  1st Qu.:36.00
Median :3.000   Median :14.00   Median :3190   Median :192.5   Median :40.00
Mean   :2.993   Mean   :13.76   Mean   :3019   Mean   :187.9   Mean   :39.65
3rd Qu.:3.500  3rd Qu.:16.75  3rd Qu.:3600  3rd Qu.:203.8  3rd Qu.:43.00
Max.   :5.000   Max.   :23.00   Max.   :4840   Max.   :233.0   Max.   :51.00
```

## 5. Kickstart

```

displacement      gear_ratio      foreign
Min.    : 79.0   Min.    :2.190   Min.    :0.0000
1st Qu.:119.0   1st Qu.:2.730   1st Qu.:0.0000
Median  :196.0   Median  :2.955   Median  :0.0000
Mean    :197.3   Mean    :3.015   Mean    :0.2973
3rd Qu.:245.2   3rd Qu.:3.352   3rd Qu.:1.0000
Max.    :425.0   Max.    :3.890   Max.    :1.0000

```

```
glimpse(auto) # Compact display of the structure of the dataset
```

```

Rows: 74
Columns: 12
$ make           <chr> "AMC Concord", "AMC Pacer", "AMC Spirit", "Buick Century"~
$ price          <dbl> 4099, 4749, 3799, 4816, 7827, 5788, 4453, 5189, 10372, 40~
$ mpg            <dbl> 22, 17, 22, 20, 15, 18, 26, 20, 16, 19, 14, 14, 21, 29, 1~
$ rep78          <dbl> 3, 3, NA, 3, 4, 3, NA, 3, 3, 3, 2, 3, 3, 4, 3, 2, 2, 3~
$ headroom        <dbl> 2.5, 3.0, 3.0, 4.5, 4.0, 4.0, 3.0, 2.0, 3.5, 3.5, 4.0, 3.~
$ trunk           <dbl> 11, 11, 12, 16, 20, 21, 10, 16, 17, 13, 20, 16, 13, 9, 20~
$ weight          <dbl> 2930, 3350, 2640, 3250, 4080, 3670, 2230, 3280, 3880, 340~
$ length          <dbl> 186, 173, 168, 196, 222, 218, 170, 200, 207, 200, 221, 20~
$ turn             <dbl> 40, 40, 35, 40, 43, 43, 34, 42, 43, 42, 44, 43, 45, 34, 4~
$ displacement    <dbl> 121, 258, 121, 196, 350, 231, 304, 196, 231, 231, 425, 35~
$ gear_ratio      <dbl> 3.58, 2.53, 3.08, 2.93, 2.41, 2.73, 2.87, 2.93, 2.93, 3.0~
$ foreign         <dbl+lbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~

```

```
print(auto, n = Inf) # Print all rows of the dataset
```

```

# A tibble: 74 x 12
  make       price   mpg rep78 headroom trunk weight length turn displacement
  <chr>     <dbl> <dbl> <dbl>   <dbl>   <dbl>   <dbl>   <dbl> <dbl>   <dbl>
1 AMC Concord 4099    22     3     2.5     11    2930    186    40     121
2 AMC Pacer   4749    17     3     3       11    3350    173    40     258
3 AMC Spirit   3799    22     NA    3       12    2640    168    35     121
4 Buick Cent~ 4816    20     3     4.5     16    3250    196    40     196
5 Buick Elec~ 7827    15     4     4       20    4080    222    43     350
6 Buick LeSa~ 5788    18     3     4       21    3670    218    43     231
7 Buick Opel   4453    26     NA    3       10    2230    170    34     304
8 Buick Regal  5189    20     3     2       16    3280    200    42     196
9 Buick Rivi~ 10372   16     3     3.5     17    3880    207    43     231
10 Buick Skyl~ 4082    19     3     3.5     13    3400    200    42     231
11 Cad. Devil~ 11385   14     3     4       20    4330    221    44     425
12 Cad. Eldor~ 14500   14     2     3.5     16    3900    204    43     350
13 Cad. Sevil~ 15906   21     3     3       13    4290    204    45     350
14 Chev. Chev~ 3299    29     3     2.5     9     2110    163    34     231
15 Chev. Impa~ 5705    16     4     4       20    3690    212    43     250
16 Chev. Mali~ 4504    22     3     3.5     17    3180    193    31     200
17 Chev. Mont~ 5104    22     2     2       16    3220    200    41     200
18 Chev. Monza 3667    24     2     2       7     2750    179    40     151

```

## 5. Kickstart

19	Chev. Nova	3955	19	3	3.5	13	3430	197	43	250
20	Dodge Colt	3984	30	5	2	8	2120	163	35	98
21	Dodge Dipl~	4010	18	2	4	17	3600	206	46	318
22	Dodge Magn~	5886	16	2	4	17	3600	206	46	318
23	Dodge St. ~	6342	17	2	4.5	21	3740	220	46	225
24	Ford Fiesta	4389	28	4	1.5	9	1800	147	33	98
25	Ford Musta~	4187	21	3	2	10	2650	179	43	140
26	Linc. Cont~	11497	12	3	3.5	22	4840	233	51	400
27	Linc. Mark~	13594	12	3	2.5	18	4720	230	48	400
28	Linc. Vers~	13466	14	3	3.5	15	3830	201	41	302
29	Merc. Bobc~	3829	22	4	3	9	2580	169	39	140
30	Merc. Coug~	5379	14	4	3.5	16	4060	221	48	302
31	Merc. Marq~	6165	15	3	3.5	23	3720	212	44	302
32	Merc. Mona~	4516	18	3	3	15	3370	198	41	250
33	Merc. XR-7	6303	14	4	3	16	4130	217	45	302
34	Merc. Zeph~	3291	20	3	3.5	17	2830	195	43	140
35	Olds 98	8814	21	4	4	20	4060	220	43	350
36	Olds Cutl ~	5172	19	3	2	16	3310	198	42	231
37	Olds Cutla~	4733	19	3	4.5	16	3300	198	42	231
38	Olds Delta~	4890	18	4	4	20	3690	218	42	231
39	Olds Omega	4181	19	3	4.5	14	3370	200	43	231
40	Olds Starf~	4195	24	1	2	10	2730	180	40	151
41	Olds Toron~	10371	16	3	3.5	17	4030	206	43	350
42	Plym. Arrow	4647	28	3	2	11	3260	170	37	156
43	Plym. Champ	4425	34	5	2.5	11	1800	157	37	86
44	Plym. Hori~	4482	25	3	4	17	2200	165	36	105
45	Plym. Sapp~	6486	26	NA	1.5	8	2520	182	38	119
46	Plym. Vola~	4060	18	2	5	16	3330	201	44	225
47	Pont. Cata~	5798	18	4	4	20	3700	214	42	231
48	Pont. Fire~	4934	18	1	1.5	7	3470	198	42	231
49	Pont. Gran~	5222	19	3	2	16	3210	201	45	231
50	Pont. Le M~	4723	19	3	3.5	17	3200	199	40	231
51	Pont. Phoe~	4424	19	NA	3.5	13	3420	203	43	231
52	Pont. Sunb~	4172	24	2	2	7	2690	179	41	151
53	Audi 5000	9690	17	5	3	15	2830	189	37	131
54	Audi Fox	6295	23	3	2.5	11	2070	174	36	97
55	BMW 320i	9735	25	4	2.5	12	2650	177	34	121
56	Datsun 200	6229	23	4	1.5	6	2370	170	35	119
57	Datsun 210	4589	35	5	2	8	2020	165	32	85
58	Datsun 510	5079	24	4	2.5	8	2280	170	34	119
59	Datsun 810	8129	21	4	2.5	8	2750	184	38	146
60	Fiat Strada	4296	21	3	2.5	16	2130	161	36	105
61	Honda Acco~	5799	25	5	3	10	2240	172	36	107
62	Honda Civic	4499	28	4	2.5	5	1760	149	34	91
63	Mazda GLC	3995	30	4	3.5	11	1980	154	33	86
64	Peugeot 604	12990	14	NA	3.5	14	3420	192	38	163
65	Renault Le~	3895	26	3	3	10	1830	142	34	79
66	Subaru	3798	35	5	2.5	11	2050	164	36	97
67	Toyota Cel~	5899	18	5	2.5	14	2410	174	36	134
68	Toyota Cor~	3748	31	5	3	9	2200	165	35	97

## 5. Kickstart

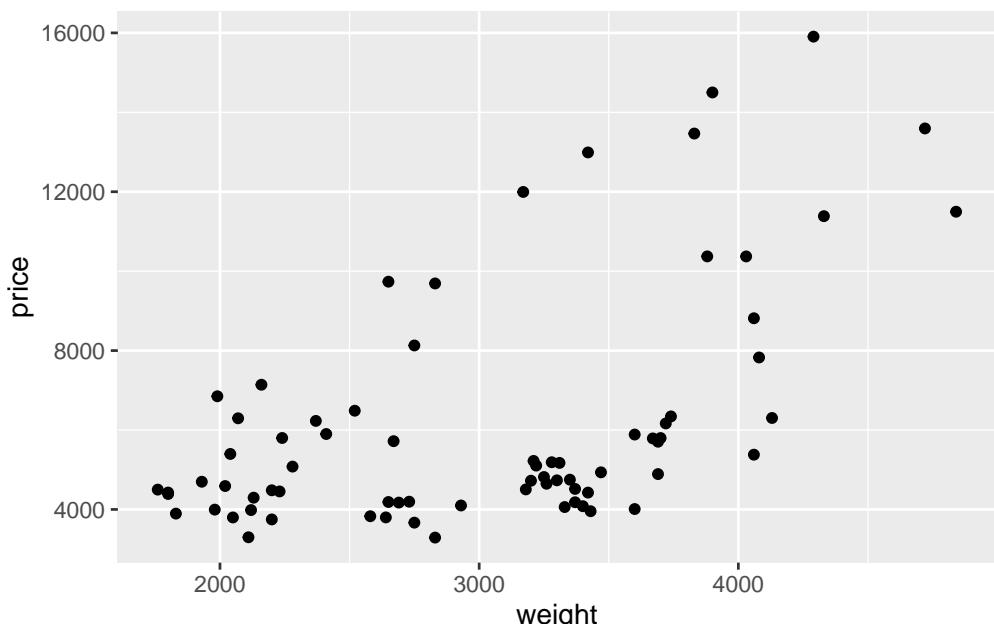
```
69 Toyota Cor~ 5719   18     5     2     11    2670    175    36      134
70 VW Dasher  7140   23     4    2.5    12    2160    172    36      97
71 VW Diesel  5397   41     5     3    15    2040    155    35      90
72 VW Rabbit  4697   25     4     3    15    1930    155    35      89
73 VW Scirocco 6850   25     4     2    16    1990    156    36      97
74 Volvo 260 11995   17     5    2.5    14    3170    193    37      163
# i 2 more variables: gear_ratio <dbl>, foreign <dbl+lbl>
```

```
# Check for duplicate entries based on the 'make' variable
auto |>
  get_dupes(make)
```

```
No duplicate combinations found of: make
```

```
# A tibble: 0 x 13
# i 13 variables: make <chr>, dupe_count <int>, price <dbl>, mpg <dbl>,
#   rep78 <dbl>, headroom <dbl>, trunk <dbl>, weight <dbl>, length <dbl>,
#   turn <dbl>, displacement <dbl>, gear_ratio <dbl>, foreign <dbl+lbl>
```

```
# Create and display a scatter plot of car price versus weight
plot_weight_price <- ggplot(auto, aes(x = weight, y = price)) +
  geom_point()
plot_weight_price
```



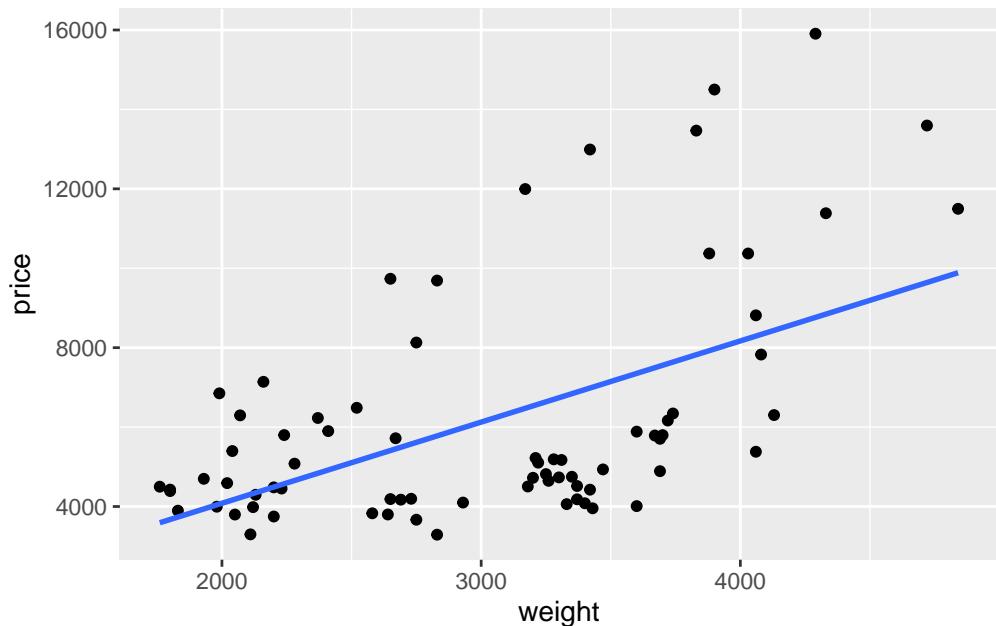
```
# Save the plot to a file
ggsave("fig/plot_weight_price.png", plot = plot_weight_price, dpi = 300)
```

```
Saving 5.5 x 3.5 in image
```

## 5. Kickstart

```
# Create a scatter plot with a linear regression line of price vs weight
plot_weight_price_fit <- ggplot(auto, aes(x = weight, y = price)) +
  geom_point() +
  geom_smooth(method = lm, se = FALSE) # 'lm' denotes linear model, 'se' is standard
plot_weight_price_fit

`geom_smooth()` using formula = 'y ~ x'
```



```
# Save the plot to a file
ggsave("fig/plot_weight_price_fit.png", plot = plot_weight_price_fit, dpi = 300)
```

Saving 5.5 x 3.5 in image

```
`geom_smooth()` using formula = 'y ~ x'
```

```
# Perform a linear regression to analyze the relationship between weight and price
reg_result <- lm(price ~ weight, data = auto)
summary(reg_result) # Display the regression results
```

Call:

```
lm(formula = price ~ weight, data = auto)
```

Residuals:

Min	1Q	Median	3Q	Max
-3341.9	-1828.3	-624.1	1232.1	7143.7

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-6.7074	1174.4296	-0.006	0.995

## 5. Kickstart

```

weight          2.0441      0.3768    5.424 7.42e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Residual standard error: 2502 on 72 degrees of freedom  
Multiple R-squared: 0.2901, Adjusted R-squared: 0.2802  
F-statistic: 29.42 on 1 and 72 DF, p-value: 7.416e-07

```

# Load and demonstrate usage of World Development Indicators (WDI) data
# Two different packages can help here:
# WDI:
# --- https://cran.r-project.org/web/packages/WDI/WDI.pdf
# --- https://github.com/vincentarelbundock/WDI
# wbstats:
# --- https://cran.r-project.org/web/packages/wbstats/wbstats.pdf
# --- https://github.com/gshs-ornl/wbstats
# ?WDI # Access documentation for the WDI package
# ?wbstats # Access documentation for the wbstats package

# Search for GDP indicators and display the first 10
WDIsearch("gdp")[1:10, ]

```

	indicator	name
689	6.0.GDP_current	GDP (current \$)
690	6.0.GDP_growth	GDP growth (annual %)
691	6.0.GDP_usd	GDP (constant 2005 \$)
692	6.0.GDPpc_constant	GDP per capita, PPP (constant 2011 international \$)
2096	BG.GSR.NFSV.GD.ZS	Trade in services (% of GDP)
2097	BG.KAC.FNEI.GD.PP.ZS	Gross private capital flows (% of GDP, PPP)
2098	BG.KAC.FNEI.GD.ZS	Gross private capital flows (% of GDP)
2099	BG.KLT.DINV.GD.PP.ZS	Gross foreign direct investment (% of GDP, PPP)
2100	BG.KLT.DINV.GD.ZS	Gross foreign direct investment (% of GDP)
2401	BI.WAG.TOTL.GD.ZS	Wage bill as a percentage of GDP

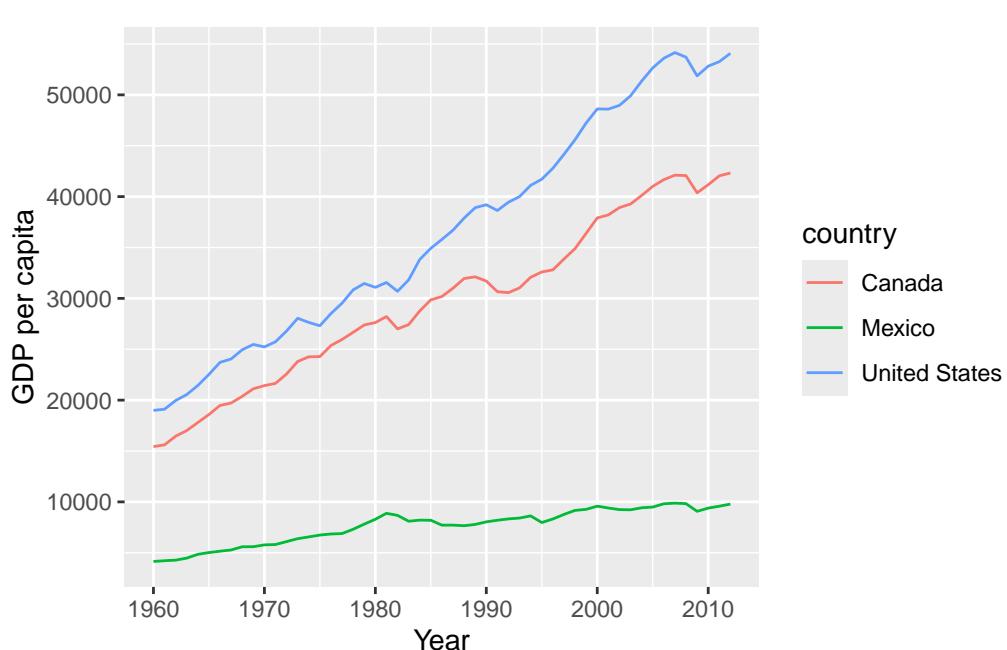
```

# Retrieve GDP per capita data for specified countries and years
df_WDI <- WDI(
  indicator = "NY.GDP.PCAP.KD",
  country = c("MX", "CA", "US"),
  start = 1960,
  end = 2012
)

# Plot GDP per capita over time for the specified countries
ggplot(df_WDI, aes(year, NY.GDP.PCAP.KD, color = country)) +
  geom_line() +
  xlab("Year") +
  ylab("GDP per capita")

```

## 5. Kickstart



```
# Retrieve GDP per capita data for specified countries and years using the wbstats package
df_wb <- wb(
  indicator = "NY.GDP.PCAP.KD",
  country = c("MX", "CA", "US"),
  start = 1960,
  end = 2012,
  return_wide = TRUE
)
```

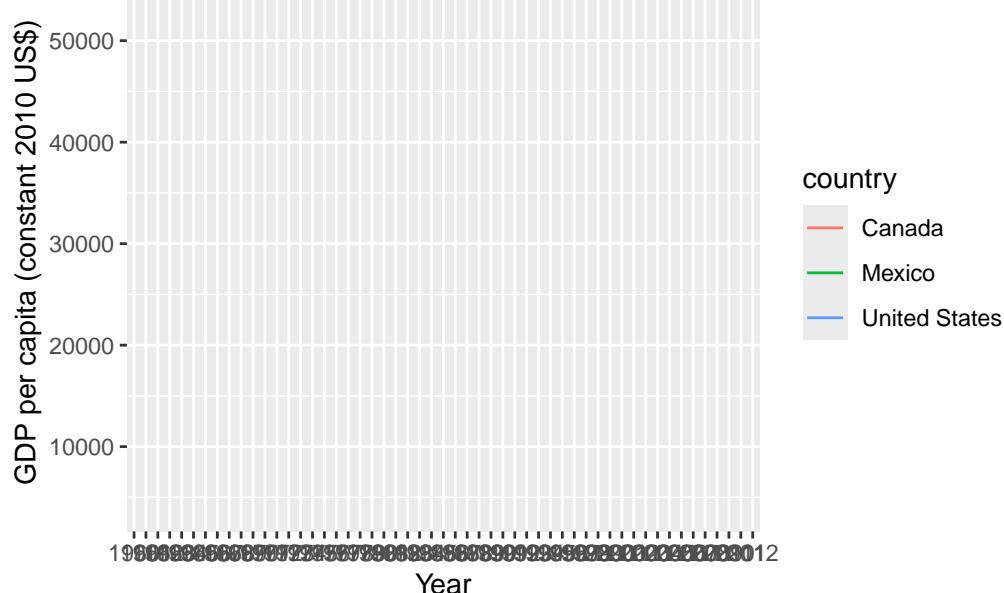
Warning: `wb()` was deprecated in wbstats 1.0.0.  
i Please use `wb\_data()` instead.

Warning: `spread\_()` was deprecated in tidyr 1.2.0.  
i Please use `spread()` instead.  
i The deprecated feature was likely used in the wbstats package.  
Please report the issue at <<https://github.com/nset-ornl/wbstats/issues>>.

```
# Plot GDP per capita over time for the specified countries
ggplot(df_wb, aes(date, NY.GDP.PCAP.KD, color = country)) +
  geom_line() +
  xlab("Year") +
  ylab("GDP per capita (constant 2010 US$)")
```

`geom\_line()`: Each group consists of only one observation.  
i Do you need to adjust the group aesthetic?

## 5. Kickstart



```
# !!! This does not work !!!  
# Why?  
  
# Look at the data types year and date are different:  
glimpse(df_WDI)
```

```
Rows: 159  
Columns: 5  
$ country      <chr> "Canada", "Canada", "Canada", "Canada", "Canada", "Can~  
$ iso2c        <chr> "CA", ~  
$ iso3c        <chr> "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", ~  
$ year         <int> 2012, 2011, 2010, 2009, 2008, 2007, 2006, 2005, 2004, 2~  
$ NY.GDP.PCAP.KD <dbl> 42320.64, 42043.64, 41164.34, 40376.42, 42067.57, 42106~
```

```
glimpse(df_wb)
```

```
Rows: 159  
Columns: 5  
$ iso3c        <chr> "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", "CAN", ~  
$ date         <chr> "1960", "1961", "1962", "1963", "1964", "1965", "1966", ~  
$ iso2c        <chr> "CA", ~  
$ country      <chr> "Canada", "Canada", "Canada", "Canada", "Canada", "Cana~  
$ NY.GDP.PCAP.KD <dbl> 15432.47, 15605.52, 16455.75, 17007.69, 17800.86, 18585~
```

## 5. Kickstart

```
# Answer: A lineplot with a character variable on the x-axis does not work!

# make the two dataframe equal:
df_wb_cln <- df_wb |>
  # Convert 'date' in df_wb from character to integer
  mutate(year = as.integer(date)) |>
  # Since 'year' has been created, remove the original 'date' column
  select(-date) |>
  # Relocate columns to organize the data frame
  relocate(country, iso2c, iso3c, year, NY.GDP.PCAP.KD)

glimpse(df_WDI)
```

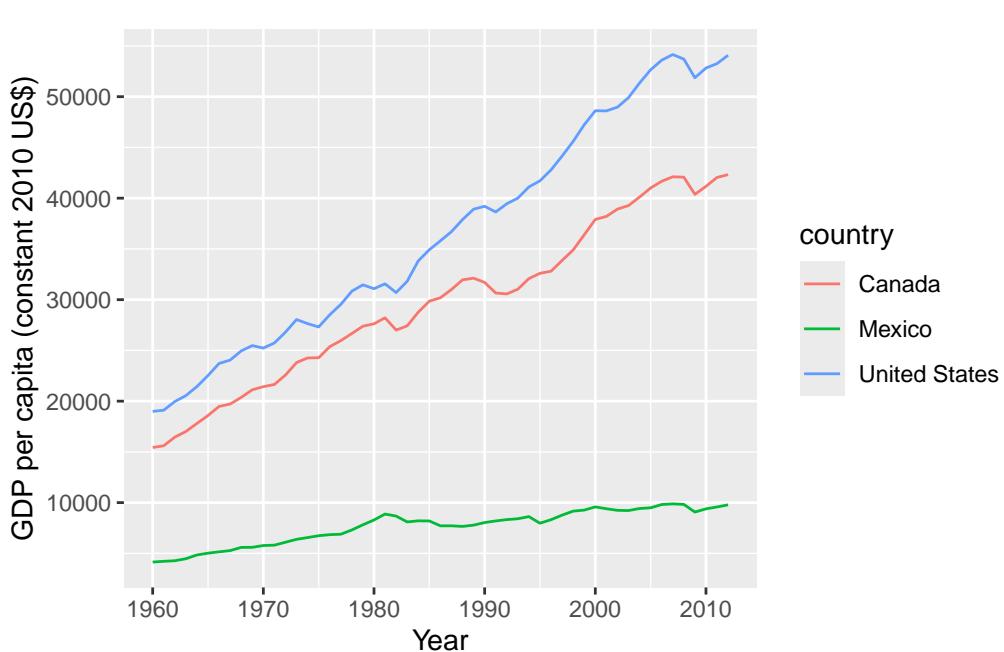
```
Rows: 159
Columns: 5
$ country      <chr> "Canada", "Canada", "Canada", "Canada", "Canada", "Can-
$ iso2c        <chr> "CA", "CA",
$ iso3c        <chr> "CAN", "CAN",
$ year         <int> 2012, 2011, 2010, 2009, 2008, 2007, 2006, 2005, 2004, 2-
$ NY.GDP.PCAP.KD <dbl> 42320.64, 42043.64, 41164.34, 40376.42, 42067.57, 42106~
```

```
glimpse(df_wb_cln)
```

```
Rows: 159
Columns: 5
$ country      <chr> "Canada", "Canada", "Canada", "Canada", "Canada", "Can-
$ iso2c        <chr> "CA", "CA",
$ iso3c        <chr> "CAN", "CAN",
$ year         <int> 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968, 1-
$ NY.GDP.PCAP.KD <dbl> 15432.47, 15605.52, 16455.75, 17007.69, 17800.86, 18585~
```

```
# Now it works:
# Plot GDP per capita over time for the specified countries
ggplot(df_wb_cln, aes(year, NY.GDP.PCAP.KD, color = country)) +
  geom_line() +
  xlab("Year") +
  ylab("GDP per capita (constant 2010 US$)")
```

## 5. Kickstart



```
suppressMessages(pacman::p_unload(  
  tidyverse, # A suite of packages designed for data science that includes tools for  
  haven, # Used for importing and exporting data with SPSS, Stata, and SAS formats.  
  janitor, # Provides functions for examining and cleaning data, such as `clean_names`  
  WDI, # Facilitates downloading data from the World Bank's World Development Indicators  
  wbstats # Provides an interface to the World Bank's APIs for a comprehensive range  
 ))
```

# 6. Pitfalls

R newbies often make the same small mistakes that can lead to major confusion, frustration and inefficiency. Some of them can be easily avoided. In this section, I will outline common pitfalls that I have repeatedly observed as an R instructor and offer practical solutions to avoid them.

## 6.1. No clue about the “working directory”

**Problem:** Students start their R sessions unaware of their current working directory. This can lead to difficulties when reading and writing files.

**Solution:** At the beginning of your R script set a working directory using `setwd()`. Consider using R Studio projects, see Section A.4. For more information, see Appendix A: *Navigating the file system* and *Workflow: scripts and projects* of Wickham and Grolemund [2023].

## 6.2. No consistent directory structure

**Problem:** Students save files in different directories without a clear scheme. This disorganization often leads to problems: Scripts and data gets lost and code breaks.

**Solution:** Organize your project into a clear directory structure from the beginning. Here is my suggestion for a directory structure but feel free to come up with your own:

Table 6.1.: Typical folder structure

Sub-Directory	What to save here
doc/	documentation
dta/	processed data
fig/	figures
lit/	literature and pdfs
ori/	original raw data that you should never change
qmd/	reports
scr/	R scripts
tab/	tables
tmp/	temporary files

This structure will save time and headaches when navigating projects.

 Tip 4: Do not save processed data unless necessary

It may seem reasonable to save data after editing, but this often isn't necessary if you're using scripts to create your data. These scripts can be rerun whenever needed, regenerating

the dataset each time. To avoid wasting disk space and maintain an organized project folder, it's advisable to save processed datasets only when the preprocessing steps are time-consuming. This way, you can keep your project folder more organized and ensure that your data analyses are always reproducible with the latest updates to your code.

### 6.3. Working manually outside R

**Problem:** Students want to get their work done quickly. This sometimes leads to them relying on manual processes that they have already mastered for their data work. This approach can lead to serious problems when it comes to the reproducibility of their data work.

Consider a typical three-step process for loading data: (1) downloading the data, (2) unpacking the data, and finally (3) importing the data into R. Many students often take a manual approach by using their Internet browser to download the data, then using their operating system's unpacking application, and finally importing the data into an R script. While this method is not inherently wrong, there is a risk that students will forget to unpack the downloaded data, resulting in them accidentally working with outdated data. In the kickstart example provided by Section 5.2, I show that all three steps can be performed seamlessly in R. This way you ensure that you are always working with the most up-to-date data.

**Solution:** Do as much as possible in the script. Invest some time to find out how to download and manipulate the data within R. If it is not possible or if alternatives are superior, describe what you do outside of R explicitly and write a warning note at the top of your script.

### 6.4. No active R Packages management

**Problem:** Students often forget to install and/or load the packages correctly at the beginning of a script. Some unnecessarily install packages repeatedly when running a script. All this can lead to errors and interruptions.

**Solution:** At the beginning of each script, make sure that all required packages are loaded correctly. Use the `pacman` package, which provides the `p_load()` function to load and, if necessary, install packages and the `p_unload(all)` function to unload all packages.

💡 Tip 5: Start your script with

```
if (!require(pacman)) install.packages("pacman")
pacman::p_unload(all)
pacman::p_load(tidyverse, janitor)
setwd("~/your-directory/")
rm(list = ls())
```

### 6.5. Confusion between console and script

**Problem:** Alternating between running code in the console and from the script without a systematic approach can lead to untracked changes and confusion about the current state of objects in the workspace. Additionally, students often borrow code snippets from others and run

## 6. Pitfalls

only the sections that seem immediately relevant. This practice can lead to errors or unexpected results, as such code often relies on previous commands or setups.

**Solution:** Develop the habit of testing small blocks of code in the console but run the complete script regularly to ensure everything works in sequence. Use shortcuts like **Ctrl + Alt + R** to source the entire script or **Ctrl + Alt + B/E** to execute it up to a specific point.

## 6.6. Misunderstanding data types and formats

**Problem:** Misusing or misunderstanding R’s data types and structures can lead to errors in data manipulation and analysis. Many functions require certain types of data. For example, the `tidyverse` packages required data to be “tidy”. Moreover, data often comes with errors and/or missings (`NA`). Beginners overlook data cleaning and considering missings.

**Solution:** Familiarize yourself with basic data types and structures like vectors, lists, data frames/tibbles, and factors. For more information, see Section 7.2 and [Data tidying of Wickham and Grolemund \[2023\]](#). Moreover, spend adequate time on data cleaning and preprocessing. Techniques such as handling missing values, normalizing data, and correcting data types are critical. For more information, see [Missing values of Wickham and Grolemund \[2023\]](#).

## 6.7. Lack of knowledge about data identification

**Problem:** Students often handle data without understanding which variables uniquely identify the information contained in other variables. It is crucial to recognize these identifying variables and verify their uniqueness to ensure data integrity.

**Solution:** Perform checks for uniqueness at the beginning of their data analysis process. See exercise *Names and duplicates* in Section 9.17 and the `get_dupes` function introduced in Section 7.4.3.2.

💡 Tip 6: Always check your data with `get_dupes`

For example, you expect that your dataframe `df` is a panel dataset. With

```
get_dupes(df, country, year)
```

you can check whether the two variables `country` and `year` indeed identify each row uniquely.

## 6.8. Losing track of data due to excessive overwriting

**Problem:** Students often manipulate their data by repeatedly overwriting the same object. This can lead to confusion about the data’s current state and the transformations applied.

**Solution:** Minimize the number of assignments to a single object. Instead, create a new object with a descriptive and concise name each time you alter the data. This practice helps maintain clarity about each stage of data manipulation.

For example, if you're working with data `df`, you might store the cleanded data as `df_cln`, then after filtering for specific criteria, you could use `df_cln_flt`, and finally, if you aggregate the data, name it `df_cln_flt_agg`. Having some clear naming convention makes it clear what each dataset represents and the transformations it has undergone.

 Tip 7: Do clear the environment at the beginning of a script with

```
rm(list = ls())
```

## 6.9. No documentation

**Problem:** Students do not comment code. This makes it hard to remember the purpose of various lines of code and difficult for other people to read and understand the code.

**Solution:** Regularly comment your code, explaining why something is done, not just what is done. Use clear, concise comments to improve readability and maintainability.

## 6.10. Ignoring error messages and warnings

**Problem:** Students see that their code doesn't work but do not read the error message which often contains hints for solving the problem.

**Solution:** Read and follow error messages. Do not ignore warnings or errors unless you know what they mean. Study what the error message might mean. Use online resources such as Google and ChatGPT, see Figure 6.1. Finally, have the confidence to implement the suggested solution. Don't be frustrated if the first attempt does not work: Try again and play around.

 Tip 8: Common error messages

Students often come to me with error messages suggesting that they install the `tinytex` package or the `RTools` compiler for Windows. Since they are not familiar with these R and software packages, they wonder if it is safe to proceed with the installation. My answer here is: yes, it is recommendable to install it.

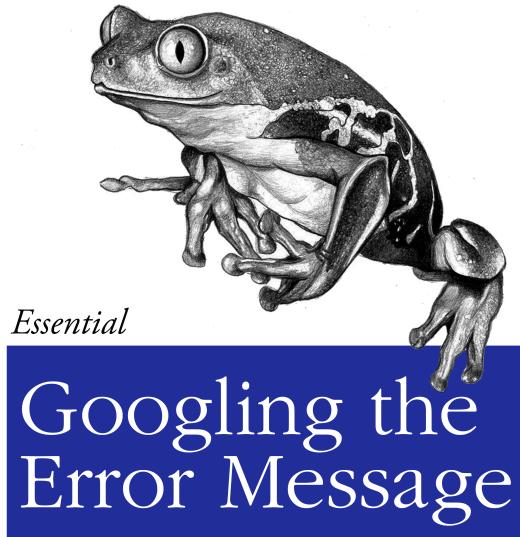
Based the [report of Noam Ross](#) who examine roughly 10,000 R error messages, the most frequently encountered error messages of R are shown in Table 6.2 with some ideas of mine what to do.

Table 6.2.: Most frequent errors

Error Type	Some suggestions on what to do
<code>Could not find function</code>	Check spelling of the function and whether the respective packages are loaded properly.
<code>Error in if</code>	This suggests an issue with non-logical or missing values in a conditional statement. Check syntax and spelling. Maybe use ChatGPT to debug the code.
<code>Error in eval</code>	Points to references to non-existent objects.
<code>Cannot open</code>	Check if the files exist at the place you try to call them.
<code>No applicable method</code>	Check your data type and whether it fits to the requirements of the functions you're trying to use.

Figure 6.1.: Googling the error message

*The internet will make those bad words go away*



O RLY?

*The Practical Developer  
@ThePracticalDev*

*Source: DEV Community on GitHub*

#### Package errors

This can stem from issues with installing, compiling, or loading a package. Maybe you try to re-install the package and their dependencies, or update the packages.

## 6.11. No attempt to identify the problem and troubleshoot

**Problem:** Students often do not fully understand the problems they encounter, which can lead to difficulties in seeking solutions. It's common for students to feel overwhelmed and seek help without attempting to find the source of the problem and without attempting to work out possible solutions first.

**Solution:** When you encounter an issue in your R code, it's crucial to methodically dissect the problem. Here's how you can effectively *troubleshoot*, that is, a problem-solving skill that is essential for becoming proficient in programming:

- **Identify the problem:** Attempt to identify the issue to better understand its nature. Once you know which line of code is causing some trouble, you are often close to a solution. Commenting out parts of your script or going back until there is no error can help here.
- **Active solution search:** Once you've identified the problem, actively look for solutions. This can include consulting the R documentation, searching for similar issues online, or asking others for help.
- **Trial and error process:** Don't hesitate to experiment with different solutions to see what works best.

## 6. Pitfalls

- **Seek Help:** If you're stuck, ask for help from more experienced R users or communities.
- **Minimal Reproducible Example (MRE):** When asking for help, explain your problem precisely and provide a MRE. That is the simplest version of the code that still produces the error, including only essential data and code. This practice not only aids in self-troubleshooting but also makes it easier for others to help by providing a clear, concise context.
- **Additional information:** Sometime the interplay of your the packages loaded, your operating system, the version of R, and/or the RStudio version may play a role in your problem. Thus, when seeking help, be sure to provide information about your machine, including the operating system, the version of R, and the packages you have loaded. You can use the `sessionInfo()` function to gather this information.

Here's an example from my machine:

```
sessionInfo()

R version 4.5.0 (2025-04-11)
Platform: x86_64-pc-linux-gnu
Running under: Debian GNU/Linux 12 (bookworm)

Matrix products: default
BLAS: /usr/lib/x86_64-linux-gnu/openblas-pthread/libblas.so.3
LAPACK: /usr/lib/x86_64-linux-gnu/openblas-pthread/libopenblas-p0.3.21.so; LAPACK version 3.21.0

locale:
[1] LC_CTYPE=en_US.UTF-8          LC_NUMERIC=C
[3] LC_TIME=en_US.UTF-8          LC_COLLATE=en_US.UTF-8
[5] LC_MONETARY=en_US.UTF-8      LC_MESSAGES=en_US.UTF-8
[7] LC_PAPER=en_US.UTF-8         LC_NAME=C
[9] LC_ADDRESS=C                 LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8   LC_IDENTIFICATION=C

time zone: Europe/Berlin
tzcode source: system (glibc)

attached base packages:
[1] stats      graphics    grDevices utils      datasets  methods   base

loaded via a namespace (and not attached):
[1] compiler_4.5.0    fastmap_1.2.0    cli_3.6.4      tools_4.5.0
[5] htmltools_0.5.8.1 rstudioapi_0.17.1 rmarkdown_2.29  knitr_1.50
[9] jsonlite_2.0.0    xfun_0.52       digest_0.6.37   rlang_1.1.6
[13] evaluate_1.0.3
```

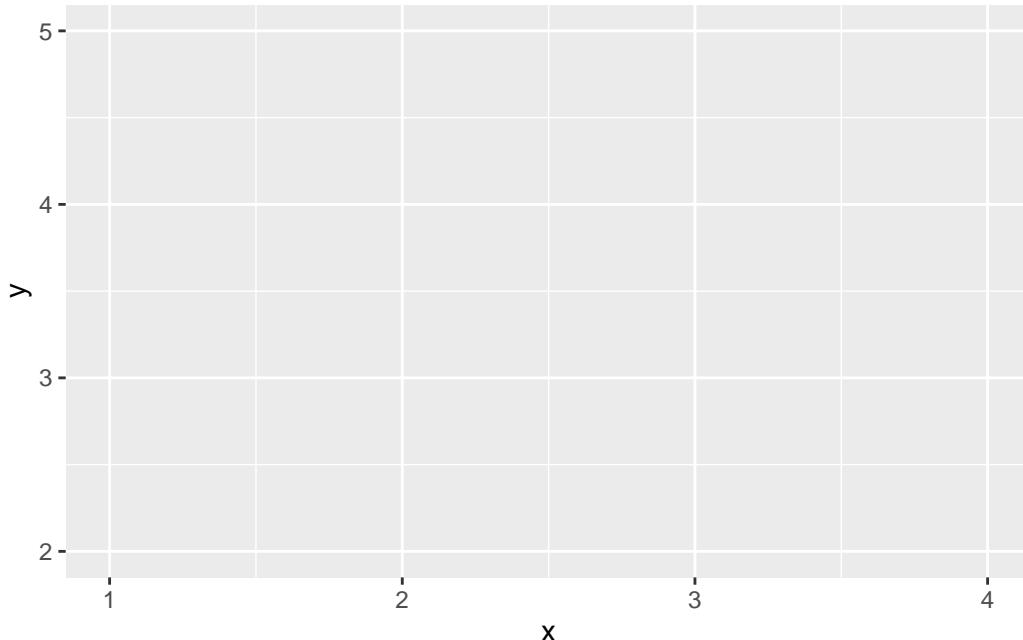
For more information, see [Workflow: getting help](#) of Wickham and Grolemund [2023].

Example of a minimal reproducible example

For example, the following script is a MRE:

## 6. Pitfalls

```
library(ggplot2)
data <- data.frame(x = 1:4, y = c(2, 3, 5, 3.4))
ggplot(data, aes(x, y))
```



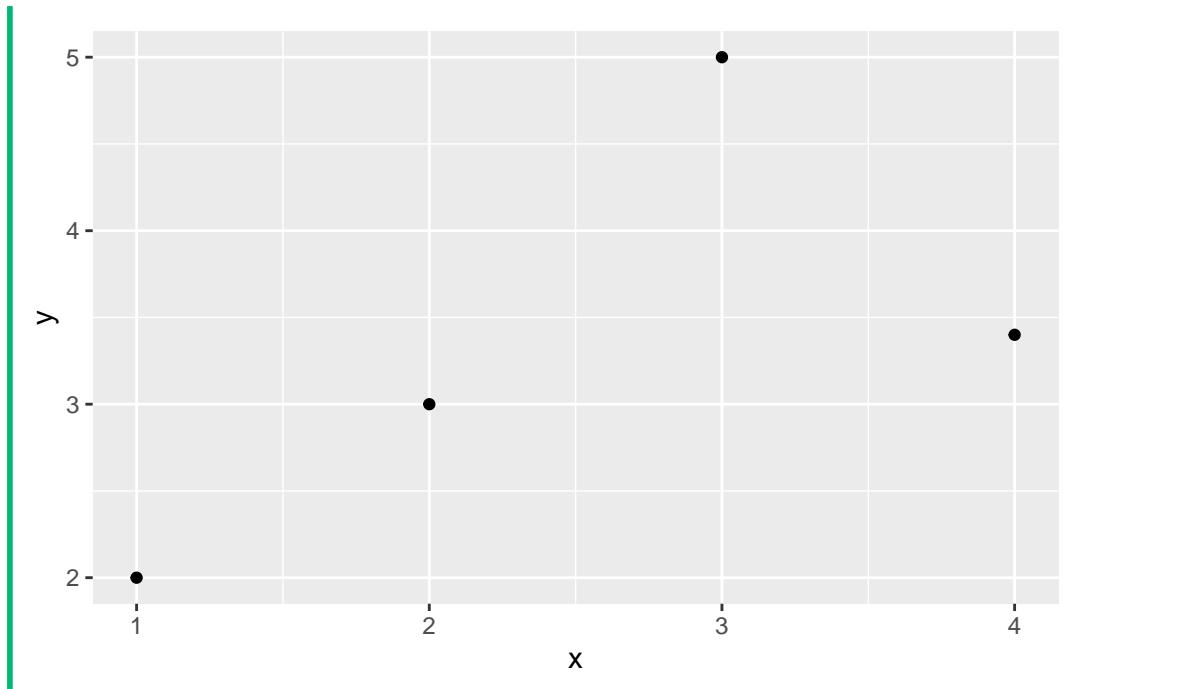
```
+geom_point()
```

Error:

```
! Cannot use `+` with a single argument.
i Did you accidentally put `+` on a new line?
```

Everybody who copies these few lines of code can reproduce the shown error message and hence can work on a solution. Obviously, the + was set falsely. It must be placed in the line of ggplot:

```
library(ggplot2)
data <- data.frame(x = 1:4, y = c(2, 3, 5, 3.4))
ggplot(data, aes(x, y)) +
  geom_point()
```



## 6.12. Unstylish code

To avoid issues while programming in R, it's essential to understand and adhere to various conventions, rules, and best practices specific to the language. Following these conventions makes your code more readable and simplifies your own experience with R. Below, you will find a non-exhaustive list of these guidelines.

1. Do remember that R programming language is case sensitive.
2. Do start names of objects such as vectors, numbers, variables, and data frames with a letter, not a number.
3. Do avoid using dots in names of objects.
4. Do avoid using certain keywords in naming objects, such as if, else, repeat, while, function, for, in, next, break, TRUE, FALSE, NULL, Inf, NaN, and NA.
5. Do use front slash / instead of backslash \ for navigating the file system (see Appendix A).
6. Do not use whitespace and indentation for naming files, directories, or objects.
7. Do define objects to represent hard-coded values instead of using them directly in code.
8. Do remember to (install and) load packages that contain functions you want to use.
9. Do use <- instead of = for assignment.

 Tip 9

There are two packages, `styler` and `lintr`, that support you writing code according to the [The tidyverse style guide](#) of Wickham [2024].

**Part III.**

**Do stuff**

# 7. Manage data

## 7.1. Import and generate data

### 7.1.1. Assigning data to an object using the assignment operator <-

Suppose I'm trying to calculate how much money I'm going to make from selling an item. Let's assume you sell 350 units. To create a variable called `sales` and assigns a value to it, we need to use the assignment operator of R, that is, `<-`:

```
sales <- 350
```

When you send that line of code to the console, it doesn't print out any output but it creates the object `sales`. In Rstudio, you can see the object in the *environment panel* at the top right. Alternatively, you can call the object in the console:

```
sales
```

```
[1] 350
```

R also allows to use `->` and `=` for the assignment. For example, the following ways of assigning data are equivalent:

```
350 -> sales
sales = 350
sales <- 350
```

However, it is common practice and “good style” to use `<-` and I recommend only to use this one because it is easier to read in scripts.

### 7.1.2. Vectors and matrices

We already got known to the `c()` function which allows to combine multiple values into a vector or list. Here are some examples how you can use this function to create vectors and matrices:

```
# defining multiple vectors using the colon operator `:`
v_a <- c(1:3)
v_a
```

```
[1] 1 2 3
```

## 7. Manage data

```
v_b <- c(10:12)
v_b
```

```
[1] 10 11 12
```

```
# creating matrix
m_ab <- matrix(c(v_a, v_b), ncol = 2)
m_cbind <- cbind(v_a, v_b)
m_rbind <- rbind(v_a, v_b)

# print matrix
print(m_ab)
```

```
 [,1] [,2]
[1,]    1   10
[2,]    2   11
[3,]    3   12
```

```
print(m_cbind)
```

```
      v_a v_b
[1,] 1 10
[2,] 2 11
[3,] 3 12
```

```
print(m_rbind)
```

```
 [,1] [,2] [,3]
v_a    1    2    3
v_b    10   11   12
```

```
# defining row names and column names
rown <- c("row_1", "row_2", "row_3")
coln <- c("col_1", "col_2")

# creating matrix
m_ab_label <- matrix(m_ab,
  ncol = 2, byrow = FALSE,
  dimnames = list(rown, coln)
)

# print matrix
print(m_ab_label)
```

```
      col_1 col_2
row_1    1   10
row_2    2   11
row_3    3   12
```

## 7. Manage data

The two most common formats to store and work with data in R are `dataframe` and `tibble`. Both formats store table-like structures of data in rows and columns. We will learn more on that in section [Section 7.2](#).

```
# convert the matrix into dataframe
df_ab <- as.data.frame(m_ab_label)
tbl_ab <- data.frame(m_ab_label)
```

### Exercise

See exercise in Section [9.1](#): *Import data with `c()`.*

### 7.1.3. Open RData files

You can save some of your objects with `save()` or all with `save.image()`. Load data that are stored in the `.RData` format can be loaded with `load()`. Please note, when you delete an object in R, you cannot recover it by clicking some *Undo button*. With `rm()` you remove objects from your workspace and with `rm(list = ls())` you clear all objects from the workspace.

### 7.1.4. Open datasets of packages

The `datasets` package contains numerous datasets that are commonly used in textbooks. To get an overview of all the datasets provided by the package, you can use the command `help(package = datasets)`. One such dataset that we will be using further is the `mtcars` dataset:

```
library("datasets")
head(mtcars, 3)

      mpg cyl disp  hp drat    wt  qsec vs am gear carb
Mazda RX4     21.0   6 160 110 3.90 2.620 16.46  0  1    4    4
Mazda RX4 Wag 21.0   6 160 110 3.90 2.875 17.02  0  1    4    4
Datsun 710    22.8   4 108  93 3.85 2.320 18.61  1  1    4    1

?mtcars # data dictionary
```

### 7.1.5. Import data using public APIs

An API which stands for *application programming interface* specifies how computers can exchange information. There are many R packages available that provide a convenient way to access data from various online sources directly within R using the API of webpages. In most cases, it's better to download and import data within R using these tools than to navigate through the website's interface. This ensures that changes can be made easily at any time and that the data is always up-to-date. For example, `wbstats` provides access to World Bank data, `eurostat` allows users to access Eurostat databases, `fredr` makes it easy to obtain data from the *Federal Reserve Economic Data (FRED)* platform, which offers economic data for the United States, `ecb` provides an interface to the European Central Bank's Statistical Data Warehouse, and the `OECD` package facilitates the extraction of data from the *Organization for Economic Cooperation and Development (OECD)*. Here is an example using the `wbstats` package:

## 7. Manage data

```
# install.packages("wbstats")
library("wbstats")
# GDP at market prices (current US$) for all available countries and regions
df_gdp <- wb(indicator = "NY.GDP.MKTP.CD")
```

Warning: `wb()` was deprecated in wbstats 1.0.0.  
i Please use `wb\_data()` instead.

```
head(df_gdp, 3)
```

```
iso3c date      value indicatorID      indicator iso2c
2   AFE 2023 1.245472e+12 NY.GDP.MKTP.CD GDP (current US$)    ZH
3   AFE 2022 1.191423e+12 NY.GDP.MKTP.CD GDP (current US$)    ZH
4   AFE 2021 1.085745e+12 NY.GDP.MKTP.CD GDP (current US$)    ZH
                           country
2 Africa Eastern and Southern
3 Africa Eastern and Southern
4 Africa Eastern and Southern
```

```
glimpse(df_gdp)
```

Rows: 14,307  
Columns: 7

```
$ iso3c      <chr> "AFE", "AFE", "AFE", "AFE", "AFE", "AFE", "AFE", "A~
$ date       <chr> "2023", "2022", "2021", "2020", "2019", "2018", "2017", "2~
$ value      <dbl> 1.245472e+12, 1.191423e+12, 1.085745e+12, 9.333918e+11, 1.~
$ indicatorID <chr> "NY.GDP.MKTP.CD", "NY.GDP.MKTP.CD", "NY.GDP.MKTP.CD", "NY.~
$ indicator    <chr> "GDP (current US$)", "GDP (current US$)", "GDP (current US~
$ iso2c      <chr> "ZH", "ZH", "ZH", "ZH", "ZH", "ZH", "ZH", "ZH", "ZH"~
$ country     <chr> "Africa Eastern and Southern", "Africa Eastern and Souther~
```

```
summary(df_gdp)
```

```
iso3c                  date          value      indicatorID
Length:14307        Length:14307      Min.   :2.586e+06  Length:14307
Class :character    Class :character   1st Qu.:2.294e+09  Class :character
Mode  :character    Mode  :character   Median :1.692e+10  Mode  :character
                           Mean   :1.185e+12
                           3rd Qu.:2.013e+11
                           Max.  :1.062e+14
indicator            iso2c          country
Length:14307        Length:14307      Length:14307
Class :character    Class :character   Class :character
Mode  :character    Mode  :character   Mode  :character
```

Figure 7.1.: The logo of the packages `readr`, `haven`, and `readxl`

### 7.1.6. Import various file formats

RStudio provides convenient data import tools that can be accessed by clicking *File > Import Dataset*. In addition, tidyverse offers packages for importing data in various formats. This [cheatsheet](#), for example, is about the packages `readr`, `readxl` and `googlesheets4`. The first allows you to read data in various file formats, including fixed-width files like `.csv` and `.tsv`. The package `readxl` can read in Excel files, i.e., `.xls` and `.xlsx` file formats and `googlesheets4` allows to read and write data from Google Sheets directly from R.

For more information, I recommend once again the second version book *R for Data Science* by [Wickham and Grolemund \[2023\]](#). In particular, check out the “[Data tidying](#)” section for importing CSV and TSV files, the “[Spreadsheets](#)” section for Excel files, the “[Databases](#)” section for retrieving data with SQL, the “[Arrow](#)” section for working with large datasets, and the “[Web scraping](#)” section for extracting data from web pages.

For an overview on packages for reading data that are provided by the tidyverse universe, see [here](#).

### 7.1.7. Examples

Flat files such as CSV (Comma-Separated Values) are among the most common and straightforward data formats to work with.

```
data_csv <- read_csv("https://github.com/hubchev/courses/raw/main/dta/classdata.csv")
```

Excel files, due to their wide use in business and research, require a specific approach specifying sheets and cell ranges.

```
BWL_Zeitschriftenliste <-
  read_excel(
    "https://www.forschungsmonitoring.org/VWL_Zeitschriftenliste%202023.xlsx",
    sheet = "SJR main",
    range = "A1:D1977"
  )
```

## 7.2. Data

### 7.2.1. Data frames and tibbles

Figure 7.2.: The logos of the `tidyr` and `tibble` packages



Both *data frames* and *tibbles* are two of the most commonly used data structures in R for handling tabular data. A tibble actually is a data frame and you can use all functions that work with a data frame also with a tibble. However, a tibble has some additional features in printing and subsetting. Please note, data frames are provided by base R while tibbles are provided by the `tidyverse` package. This means that if you want to use tibbles you must load `tidyverse`. It turned out that it is helpful that a tibble has the following features to simplify working with data:

- Each vector is labeled by the variable name.
- Variable names don't have spaces and are not put in quotes.
- All variables have the same length.
- Each variable is of a single type (numeric, character, logical, or a categorical).

### 7.2.2. Tidy data

A popular quote from Hadley Wickham is that

“tidy datasets are all alike, but every messy dataset is messy in its own way” [Hadley, 2014, p. 2].

It paraphrases the fact that it is a good idea to set rules how a dataset should structure its information to make it easier to work with the data. The tidyverse requires the data to be structured like is illustrated in Figure Figure 7.3. The rules are:

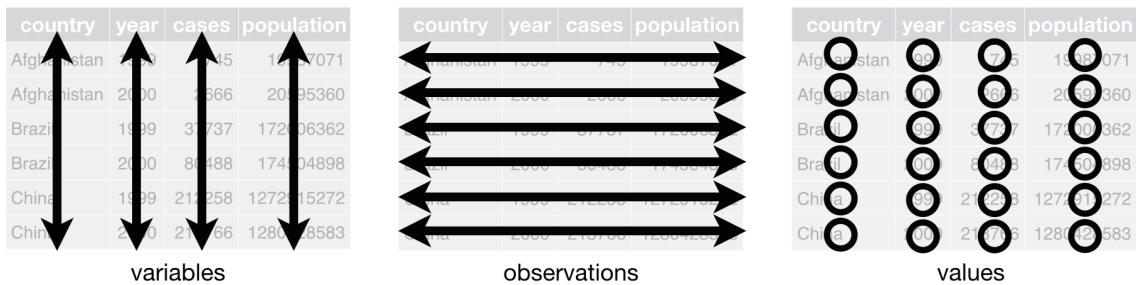
1. Each variable is a column and vice versa.
2. Each observation is a row and vice versa.
3. Each value is a cell.

Whenever data follow that consistent structure, we speak of *tidy data*. The underlying uniformity of tidy data facilitates learning and using data manipulation tools.

One difference between data frames and tibbles is that dataframes store the row names. For example, take the `mtcars` dataset which consists of 32 different cars and the names of the cars are not stored as rownames:

## 7. Manage data

Figure 7.3.: Features of a tidy dataset: variables are columns, observations are rows, and values are cells



Source: *Wickham and Grolemund [2023]*.

```
class(mtcars) # mtcars is a data frame
```

```
[1] "data.frame"
```

```
rownames(mtcars)
```

```
[1] "Mazda RX4"           "Mazda RX4 Wag"      "Datsun 710"
[4] "Hornet 4 Drive"      "Hornet Sportabout" "Valiant"
[7] "Duster 360"          "Merc 240D"         "Merc 230"
[10] "Merc 280"            "Merc 280C"         "Merc 450SE"
[13] "Merc 450SL"          "Merc 450SLC"       "Cadillac Fleetwood"
[16] "Lincoln Continental" "Chrysler Imperial" "Fiat 128"
[19] "Honda Civic"         "Toyota Corolla"    "Toyota Corona"
[22] "Dodge Challenger"   "AMC Javelin"      "Camaro Z28"
[25] "Pontiac Firebird"   "Fiat X1-9"        "Porsche 914-2"
[28] "Lotus Europa"        "Ford Pantera L"   "Ferrari Dino"
[31] "Maserati Bora"       "Volvo 142E"        "
```

To store `mtcars` as a tibble, we can use the `as_tibble` function:

```
tbl_mtcars <- as_tibble(mtcars)
class(tbl_mtcars) # check if it is a tibble now
```

```
[1] "tbl_df"     "tbl"        "data.frame"
```

```
is_tibble(tbl_mtcars) # alternative check
```

```
[1] TRUE
```

```
head(tbl_mtcars, 3)
```

## 7. Manage data

```
# A tibble: 3 x 11
  mpg   cyl  disp    hp  drat    wt  qsec    vs    am  gear  carb
  <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 21      6    160   110   3.9   2.62  16.5     0     1     4     4
2 21      6    160   110   3.9   2.88  17.0     0     1     4     4
3 22.8    4    108    93   3.85  2.32  18.6     1     1     4     1
```

When we look at the data, we've lost the names of the cars. To store these, you need to first add a column to the dataframe containing the rownames and then you can generate the tibble:

```
tbl_mtcars <- mtcars |>
  rownames_to_column(var = "car") |>
  as_tibble()
class(tbl_mtcars)
```

```
[1] "tbl_df"     "tbl"        "data.frame"
```

```
head(tbl_mtcars, 3)
```

```
# A tibble: 3 x 12
  car       mpg   cyl  disp    hp  drat    wt  qsec    vs    am  gear  carb
  <chr>     <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 Mazda RX4  21      6    160   110   3.9   2.62  16.5     0     1     4     4
2 Mazda RX4 W~ 21      6    160   110   3.9   2.88  17.0     0     1     4     4
3 Datsun 710  22.8    4    108    93   3.85  2.32  18.6     1     1     4     1
```

### 7.2.3. Data types

In R, different data classes, or types of data exist:

- *numeric*: can be any real number
- *character*: strings and characters
- *integer*: any whole numbers
- *factor*: any categorical or qualitative variable with finite number of distinct outcomes
- *logical*: contain either TRUE or FALSE
- *Date*: special format that describes time

The following example should exemplify these types of data:

```
integer_var <- c(1, 2, 3, 4, 5)
numeric_var <- c(1.1, 2.2, NA, 4.4, 5.5)
character_var <- c("apple", "banana", "orange", "cherry", "grape")
factor_var <- factor(c("red", "yellow", "red", "blue", "green"))
logical_var <- c(TRUE, TRUE, TRUE, FALSE, TRUE)
date_var <- as.Date(c("2022-01-01", "2022-02-01", "2022-03-01", "2022-04-01", "2022-05-01"))

date_var[2] - date_var[5] # number of days in between these two dates
```

Time difference of -89 days

There are some special data values used in R that needs further explanation:

- `NA` stands for *not available* or *missing* and is used to represent missing or undefined values.
- `Inf` stands for *infinity* and is used to represent mathematical infinity, such as the result of dividing a non-zero number by zero. Can be positive or negative.
- `NULL` represents an empty or non-existent object. It is often used as a placeholder when a value or object is not yet available or when an object is intentionally removed.
- `NaN` stands for *not a number* and is used to represent an undefined or unrepresentable value, such as the result of taking the square root of a negative number. It can also occur as a result of certain arithmetic operations that are undefined. In contrast to `NA` it can only exist in numerical data.

## 7.3. Operators

An overview of the most important operators of R is provided in Appendix B.

### 7.3.1. Algebraic operators

R can perform any kind of arithmetic calculation using the operators listed in Table 7.1.

Table 7.1.: Basic algebraic operators

Operation	Operator	Example input	Example output
addition	<code>+</code>	<code>10+2</code>	12
subtraction	<code>-</code>	<code>9-3</code>	6
multiplication	<code>*</code>	<code>5*5</code>	25
division	<code>/</code>	<code>10/3</code>	3
power	<code>^</code>	<code>5^2</code>	25

### 7.3.2. The pipe operator: `|>`

The pipe operator, `%>%`, comes from the `magrittr` package, which is also part of the tidyverse package. The pipe operator, `|>`, has been part of base R since version 4.1.0. For most cases, these two operators are identical. The pipe operator is designed to help you write code in a way that is easier to read and understand. As R is a functional language, code often contains a lot of parentheses, ( and ). Nesting these parentheses together can be complex and make your R code hard to read and understand, which is where `|>` comes to the rescue! It allows you to use the output of a function as the input of the next function.

 Tip 10: Set the native pipe in RStudio

With the keyboard shortcut `Ctrl+Shift+M`, RStudio inserts `%>%`. To change that behavior, simply check the box labeled “Use native pipe operator, `|>`” in the Global Options, see: `Tools > Global Options > Code > Editing`.

Consider the following example of code to explain the usage of the pipe operator:

## 7. Manage data

```
# create some data `x`
x <- c(1, 1.002, 1.004, .99, .99)
# take the logarithm of `x`,
log_x <- log(x)
# compute the lagged and iterated differences (see `diff()`)
growth_rate_x <- diff(log_x)
growth_rate_x
```

```
[1] 0.001998003 0.001994019 -0.014042357 0.000000000
```

```
# round the result (4 digit)
growth_rate_x_round <- round(growth_rate_x, 4)
growth_rate_x_round
```

```
[1] 0.002 0.002 -0.014 0.000
```

That is rather long and we actually don't need objects `log_x`, `growth_rate_x`, and `growth_rate_x_round`. Well, then let us write that in a nested function:

```
round(diff(log(x)), 4)
```

```
[1] 0.002 0.002 -0.014 0.000
```

This is short but hard to read and understand. The solution is the “pipe”:

```
# load one of these packages: `magrittr` or `tidyverse`
library(tidyverse)

# Perform the same computations on `x` as above
x |>
  log() |>
  diff() |>
  round(4)
```

```
[1] 0.002 0.002 -0.014 0.000
```

You can read the `|>` with “*and then*” because it takes the results of some function “*and then*” does something with that in the next. For example, reading out loud the following code would sound something like this:

- I take the mtcars data, *and then*
- I consider only cars with more than 4 cylinders, *and then*
- I group the cars by the number of cylinders the cars have, *and then*
- I summarize the data and show the means of miles per gallon (mpg) and horse powers (hp) by groups of cars that distinguish by their number of cylinders.

```
mtcars |>
  filter(cyl > 4) |>
  group_by(cyl) |>
  summarise_at(c("mpg", "hp"), mean)
```

```
# A tibble: 2 x 3
  cyl   mpg    hp
  <dbl> <dbl> <dbl>
1     6 19.7 122.
2     8 15.1 209.
```

### Exercise

See exercise in Section 9.3: *Base R, %in% operator, and the pipe />*.

### 7.3.3. The %in% operator

%in% is used to subset a vector by comparison. Here's an example:

```
x <- c(1, 3, 5, 7)
y <- c(2, 4, 6, 8)
z <- c(1, 2, 3)

x %in% y
```

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

```
x %in% z
```

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

```
z %in% x
```

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

The %in% operator can be used in combination with other functions like `subset()` and `filter()`.

### Exercise

See exercise in Section 9.3: *Base R, %in% operator, and the pipe />*.

### 7.3.4. Extract operators

The *extract operators* are used to retrieve data from objects in R. The operator may take four forms, including `[]`, `[[]]`, and `$`.

`[]` allows to extract content from vector, lists, or data frames. For example,

## 7. Manage data

```
a <- mtcars[3, ]  
b <- mtcars["Datsun 710", ]  
identical(a, b)
```

```
[1] TRUE
```

```
a
```

```
  mpg cyl disp hp drat   wt  qsec vs am gear carb  
Datsun 710 22.8     4 108 93 3.85 2.32 18.61  1  1     4     1
```

extracts the third observation of the mtcars dataset, and

```
c <- mtcars[, "cyl"]  
d <- mtcars[, 2]  
identical(x, y)
```

```
[1] FALSE
```

```
c
```

```
[1] 6 6 4 6 8 6 8 4 4 6 6 8 8 8 8 8 4 4 4 4 8 8 8 8 4 4 4 8 6 8 4
```

extracts the variable/vector cyl.

The operators, [ ] and \$ extract a single item from an object. It is used to refer to an element in a list or a column in a data frame. For example,

```
e <- mtcars$cyl  
f <- mtcars[["cyl"]]  
identical(e, f)
```

```
[1] TRUE
```

```
e
```

```
[1] 6 6 4 6 8 6 8 4 4 6 6 8 8 8 8 8 4 4 4 4 8 8 8 8 4 4 4 8 6 8 4
```

will return the values of the variable cyl from the data frame mtcars. Thus, x\$y is actually just a short form for x[["y"]].

### 7.3.5. Logical operators

The extract operators can be combined with the *logical operators* (more precisely, I should call these *binary relational operators*) that are shown in Table 7.2.

Table 7.2.: Logical operators

operation	operator	example input	answer
less than	<	2 < 3	TRUE
less than or equal to	<=	2 <= 2	TRUE
greater than	>	2 > 3	FALSE
greater than or equal to	>=	2 >= 2	TRUE
equal to	==	2 == 3	FALSE
not equal to	!=	2 != 3	TRUE
not	!	!(1==1)	FALSE
or		(1==1)   (2==3)	TRUE
and	&	(1==1) & (2==3)	FALSE

Here are some examples: Select rows where the number of cylinders is greater than or equal to 6:

```
mtcars[mtcars$cyl >= 6, ]
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4
Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1
Hornet Sportabout	18.7	8	360.0	175	3.15	3.440	17.02	0	0	3	2
Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3	1
Duster 360	14.3	8	360.0	245	3.21	3.570	15.84	0	0	3	4
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4	4
Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	0	4	4
Merc 450SE	16.4	8	275.8	180	3.07	4.070	17.40	0	0	3	3
Merc 450SL	17.3	8	275.8	180	3.07	3.730	17.60	0	0	3	3
Merc 450SLC	15.2	8	275.8	180	3.07	3.780	18.00	0	0	3	3
Cadillac Fleetwood	10.4	8	472.0	205	2.93	5.250	17.98	0	0	3	4
Lincoln Continental	10.4	8	460.0	215	3.00	5.424	17.82	0	0	3	4
Chrysler Imperial	14.7	8	440.0	230	3.23	5.345	17.42	0	0	3	4
Dodge Challenger	15.5	8	318.0	150	2.76	3.520	16.87	0	0	3	2
AMC Javelin	15.2	8	304.0	150	3.15	3.435	17.30	0	0	3	2
Camaro Z28	13.3	8	350.0	245	3.73	3.840	15.41	0	0	3	4
Pontiac Firebird	19.2	8	400.0	175	3.08	3.845	17.05	0	0	3	2
Ford Pantera L	15.8	8	351.0	264	4.22	3.170	14.50	0	1	5	4
Ferrari Dino	19.7	6	145.0	175	3.62	2.770	15.50	0	1	5	6
Maserati Bora	15.0	8	301.0	335	3.54	3.570	14.60	0	1	5	8

Select rows where the number of cylinders is either 4 or 6:

## 7. Manage data

```
mtcars[mtcars$cyl == 4 | mtcars$cyl == 6, ]
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1
Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3	1
Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4	4
Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	0	4	4
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	1	4	1
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4	2
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	1	4	1
Toyota Corona	21.5	4	120.1	97	3.70	2.465	20.01	1	0	3	1
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	1	4	1
Porsche 914-2	26.0	4	120.3	91	4.43	2.140	16.70	0	1	5	2
Lotus Europa	30.4	4	95.1	113	3.77	1.513	16.90	1	1	5	2
Ferrari Dino	19.7	6	145.0	175	3.62	2.770	15.50	0	1	5	6
Volvo 142E	21.4	4	121.0	109	4.11	2.780	18.60	1	1	4	2

Select rows where the number of cylinders is 4 and the mpg is greater than 22:

```
mtcars[mtcars$cyl == 4 & mtcars$mpg > 22, ]
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	1	4	1
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4	2
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	1	4	1
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	1	4	1
Porsche 914-2	26.0	4	120.3	91	4.43	2.140	16.70	0	1	5	2
Lotus Europa	30.4	4	95.1	113	3.77	1.513	16.90	1	1	5	2

Select rows where the weight is less than 3.5 or the number of gears is greater than 4:

```
mtcars[mtcars$wt < 3.5 | mtcars$gear > 4, ]
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1
Hornet Sportabout	18.7	8	360.0	175	3.15	3.440	17.02	0	0	3	2
Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3	1

## 7. Manage data

Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4	4
Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	0	4	4
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	1	4	1
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4	2
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	1	4	1
Toyota Corona	21.5	4	120.1	97	3.70	2.465	20.01	1	0	3	1
AMC Javelin	15.2	8	304.0	150	3.15	3.435	17.30	0	0	3	2
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	1	4	1
Porsche 914-2	26.0	4	120.3	91	4.43	2.140	16.70	0	1	5	2
Lotus Europa	30.4	4	95.1	113	3.77	1.513	16.90	1	1	5	2
Ford Pantera L	15.8	8	351.0	264	4.22	3.170	14.50	0	1	5	4
Ferrari Dino	19.7	6	145.0	175	3.62	2.770	15.50	0	1	5	6
Maserati Bora	15.0	8	301.0	335	3.54	3.570	14.60	0	1	5	8
Volvo 142E	21.4	4	121.0	109	4.11	2.780	18.60	1	1	4	2

Select rows where either mpg is greater than 25 or carb is less than 2, and the number of cylinders is either 4 or 8.

```
mtcars[(mtcars$mpg > 25 | mtcars$carb < 2) & mtcars$cyl %in% c(4, 8), ]
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	1	4	1
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4	2
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	1	4	1
Toyota Corona	21.5	4	120.1	97	3.70	2.465	20.01	1	0	3	1
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	1	4	1
Porsche 914-2	26.0	4	120.3	91	4.43	2.140	16.70	0	1	5	2
Lotus Europa	30.4	4	95.1	113	3.77	1.513	16.90	1	1	5	2

## 7.4. Data manipulation

### 7.4.1. dplyr: A human readable grammar of data manipulation

Figure 7.4.: The logo of the dplyr package



The dplyr package is part of tidyverse and makes data manipulation easy as it works well with the pipe operator `|>`. The most important function are the following:

## 7. Manage data

- Reorder the rows with `arrange()`.
- Pick observations by their values with `filter()`.
- Pick variables by their names with `select()`.
- Create new variables with functions of existing variables with `mutate()`.
- Collapse many values down to a single summary with `summarise()`.
- Rename variables with `rename()`.
- Change the position of variables with `relocate()`.

These functions can be used in conjunction with `group_by()` and/or `rowwise()`, which changes the scope of each function from operating on the entire dataset to operating on it group-by-group or by rows. Moreover, you can check for conditions and take action with, for example, `if_else()` and `case_when()`.

All functions work similarly:

1. The first argument is a data frame.
2. The subsequent arguments describe what to do with the data frame.
3. The result is a new data frame.

Read the vignette of `dplyr` that you find [here](#) or with:

```
vignette("dplyr")
```

Here are some examples that may help to understand these functions:

```
library(tidyverse)

# load mtcars dataset
data(mtcars)

# filter only cars with four gears
mtcars_gear_4 <- mtcars |>
  filter(gear == 4)

# arrange rows by mpg in descending order
mtcars_gear_4 |>
  arrange(desc(mpg))
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	1	4	1
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	1	4	1
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4	2
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	1	4	1
Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2
Volvo 142E	21.4	4	121.0	109	4.11	2.780	18.60	1	1	4	2
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4	4
Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	0	4	4

## 7. Manage data

```
# Change the order of the variables
glimpse(mtcars_gear_4)
```

```
Rows: 12
Columns: 11
$ mpg <dbl> 21.0, 21.0, 22.8, 24.4, 22.8, 19.2, 17.8, 32.4, 30.4, 33.9, 27.3,~
$ cyl <dbl> 6, 6, 4, 4, 4, 6, 6, 4, 4, 4, 4
$ disp <dbl> 160.0, 160.0, 108.0, 146.7, 140.8, 167.6, 167.6, 78.7, 75.7, 71.1~
$ hp <dbl> 110, 110, 93, 62, 95, 123, 123, 66, 52, 65, 66, 109
$ drat <dbl> 3.90, 3.90, 3.85, 3.69, 3.92, 3.92, 3.92, 4.08, 4.93, 4.22, 4.08,~
$ wt <dbl> 2.620, 2.875, 2.320, 3.190, 3.150, 3.440, 3.440, 2.200, 1.615, 1.~
$ qsec <dbl> 16.46, 17.02, 18.61, 20.00, 22.90, 18.30, 18.90, 19.47, 18.52, 19~
$ vs <dbl> 0, 0, 1, 1, 1, 1, 1, 1, 1, 1
$ am <dbl> 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 1
$ gear <dbl> 4, 4, 4, 4, 4, 4, 4, 4, 4, 4
$ carb <dbl> 4, 4, 1, 2, 2, 4, 4, 1, 2, 1, 1, 2
```

```
mtcars_gear_4 |>
  relocate(cyl, disp, carb) |>
  glimpse()
```

```
Rows: 12
Columns: 11
$ cyl <dbl> 6, 6, 4, 4, 4, 6, 6, 4, 4, 4, 4
$ disp <dbl> 160.0, 160.0, 108.0, 146.7, 140.8, 167.6, 167.6, 78.7, 75.7, 71.1~
$ carb <dbl> 4, 4, 1, 2, 2, 4, 4, 1, 2, 1, 1, 2
$ mpg <dbl> 21.0, 21.0, 22.8, 24.4, 22.8, 19.2, 17.8, 32.4, 30.4, 33.9, 27.3,~
$ hp <dbl> 110, 110, 93, 62, 95, 123, 123, 66, 52, 65, 66, 109
$ drat <dbl> 3.90, 3.90, 3.85, 3.69, 3.92, 3.92, 3.92, 4.08, 4.93, 4.22, 4.08,~
$ wt <dbl> 2.620, 2.875, 2.320, 3.190, 3.150, 3.440, 3.440, 2.200, 1.615, 1.~
$ qsec <dbl> 16.46, 17.02, 18.61, 20.00, 22.90, 18.30, 18.90, 19.47, 18.52, 19~
$ vs <dbl> 0, 0, 1, 1, 1, 1, 1, 1, 1, 1
$ am <dbl> 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 1
$ gear <dbl> 4, 4, 4, 4, 4, 4, 4, 4, 4, 4
```

```
mtcars_gear_4 |>
  relocate(sort(names(mtcars_gear_4))) |>
  glimpse()
```

```
Rows: 12
Columns: 11
$ am <dbl> 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1
$ carb <dbl> 4, 4, 1, 2, 2, 4, 4, 1, 2, 1, 1, 2
$ cyl <dbl> 6, 6, 4, 4, 4, 6, 6, 4, 4, 4, 4, 4
$ disp <dbl> 160.0, 160.0, 108.0, 146.7, 140.8, 167.6, 167.6, 78.7, 75.7, 71.1~
$ drat <dbl> 3.90, 3.90, 3.85, 3.69, 3.92, 3.92, 3.92, 4.08, 4.93, 4.22, 4.08,~
$ gear <dbl> 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4
$ hp <dbl> 110, 110, 93, 62, 95, 123, 123, 66, 52, 65, 66, 109
$ mpg <dbl> 21.0, 21.0, 22.8, 24.4, 22.8, 19.2, 17.8, 32.4, 30.4, 33.9, 27.3,~
```

## 7. Manage data

```
$ qsec <dbl> 16.46, 17.02, 18.61, 20.00, 22.90, 18.30, 18.90, 19.47, 18.52, 19~  
$ vs    <dbl> 0, 0, 1, 1, 1, 1, 1, 1, 1, 1  
$ wt    <dbl> 2.620, 2.875, 2.320, 3.190, 3.150, 3.440, 3.440, 2.200, 1.615, 1.~
```

```
# filter rows where cyl = 4  
mtcars_gear_4 |>  
filter(cyl == 4)
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	1	4	1
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4	2
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	1	4	1
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	1	4	1
Volvo 142E	21.4	4	121.0	109	4.11	2.780	18.60	1	1	4	2

```
# select columns mpg, cyl, and hp  
mtcars_gear_4 |>  
select(mpg, cyl, hp) |>  
head()
```

	mpg	cyl	hp
Mazda RX4	21.0	6	110
Mazda RX4 Wag	21.0	6	110
Datsun 710	22.8	4	93
Merc 240D	24.4	4	62
Merc 230	22.8	4	95
Merc 280	19.2	6	123

```
# select columns all variables except wt and hp  
mtcars_gear_4 |>  
select(-wt, -hp) |>  
head()
```

	mpg	cyl	disp	drat	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160.0	3.90	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160.0	3.90	17.02	0	1	4	4
Datsun 710	22.8	4	108.0	3.85	18.61	1	1	4	1
Merc 240D	24.4	4	146.7	3.69	20.00	1	0	4	2
Merc 230	22.8	4	140.8	3.92	22.90	1	0	4	2
Merc 280	19.2	6	167.6	3.92	18.30	1	0	4	4

```
# select only variables starting with `c`  
mtcars_gear_4 |>  
select(starts_with("c"))
```

## 7. Manage data

```
cyl carb
Mazda RX4      6   4
Mazda RX4 Wag  6   4
Datsun 710     4   1
Merc 240D      4   2
Merc 230       4   2
Merc 280       6   4
Merc 280C      6   4
Fiat 128       4   1
Honda Civic    4   2
Toyota Corolla 4   1
Fiat X1-9      4   1
Volvo 142E     4   2
```

```
# summarize avg mpg by number of cylinders
mtcars_gear_4 |>
  group_by(cyl) |>
  summarize(avg_mpg = mean(mpg))
```

```
# A tibble: 2 x 2
  cyl avg_mpg
  <dbl>   <dbl>
1     4     26.9
2     6     19.8
```

```
# create new column wt_kg, which is wt in kg
mtcars_gear_4 |>
  select(wt) |>
  mutate(wt_kg = wt / 2.205) |>
  head()
```

```
      wt      wt_kg
Mazda RX4     2.620 1.188209
Mazda RX4 Wag 2.875 1.303855
Datsun 710    2.320 1.052154
Merc 240D     3.190 1.446712
Merc 230      3.150 1.428571
Merc 280      3.440 1.560091
```

```
# Create a new variable by calculating hp divided by wt
mtcars_new <- mtcars |>
  select(wt, hp) |>
  mutate(hp_per_t = hp / wt) |>
  head()
```

```
# Print the first few rows of the updated dataset
head(mtcars_new)
```

```
      wt      hp hp_per_t
```

## 7. Manage data

```
Mazda RX4           2.620 110 41.98473
Mazda RX4 Wag      2.875 110 38.26087
Datsun 710          2.320  93 40.08621
Hornet 4 Drive      3.215 110 34.21462
Hornet Sportabout   3.440 175 50.87209
Valiant             3.460 105 30.34682
```

```
# Rename hp to horsepower
mtcars_gear_4 |>
  rename(horsepower = hp) |>
  glimpse()
```

```
Rows: 12
Columns: 11
$ mpg            <dbl> 21.0, 21.0, 22.8, 24.4, 22.8, 19.2, 17.8, 32.4, 30.4, 33.9, ~
$ cyl             <dbl> 6, 6, 4, 4, 4, 6, 6, 4, 4, 4, 4, 4
$ disp            <dbl> 160.0, 160.0, 108.0, 146.7, 140.8, 167.6, 167.6, 78.7, 75.7 ~
$ horsepower      <dbl> 110, 110, 93, 62, 95, 123, 123, 66, 52, 65, 66, 109
$ drat            <dbl> 3.90, 3.90, 3.85, 3.69, 3.92, 3.92, 3.92, 4.08, 4.93, 4.22, ~
$ wt              <dbl> 2.620, 2.875, 2.320, 3.190, 3.150, 3.440, 3.440, 2.200, 1.6 ~
$ qsec            <dbl> 16.46, 17.02, 18.61, 20.00, 22.90, 18.30, 18.90, 19.47, 18. ~
$ vs              <dbl> 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1
$ am              <dbl> 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 1
$ gear            <dbl> 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4
$ carb            <dbl> 4, 4, 1, 2, 2, 4, 4, 1, 2, 1, 1, 2
```

### Exercise

See exercise:

- Section 9.4: *Generate and drop variables*
- Section 9.5: *Subsetting*

### 7.4.2. If statements

In many cases, it's necessary to execute certain code only when a particular condition is met. To achieve this, there are several conditional statements that can be used in code. These include:

- The `if` statement: This is used to execute a block of code if a specified condition is true.
- The `else` statement: This is used to execute a block of code if the same condition is false.
- The `else if` statement: This is used to specify a new condition to test if the first condition is false.
- The `if_else()` function: This is used to check a condition for every element of a vector.

The following examples should exemplify how these statements work:

```
# Example of if statement
if (mean(mtcars$mpg) > 20) {
  print("The average miles per gallon is greater than 20.")
}
```

## 7. Manage data

```
[1] "The average miles per gallon is greater than 20."
```

```
# Example of if-else statement
if (mean(mtcars$mpg) > 20) {
  print("The average miles per gallon is greater than 20.")
} else {
  print("The average miles per gallon is less than or equal to 20.")
}
```

```
[1] "The average miles per gallon is greater than 20."
```

```
# Example of if-else if statement
if (mean(mtcars$mpg) > 25) {
  print("The average miles per gallon is greater than 25.")
} else if (mean(mtcars$mpg) > 20) {
  print("The average miles per gallon is between 20 and 25.")
} else {
  print("The average miles per gallon is less than or equal to 20.")
}
```

```
[1] "The average miles per gallon is between 20 and 25."
```

```
# Example of if_else function
mtcars_2 <- mtcars
mtcars_2$mpg_category <- if_else(mtcars_2$mpg > 20, "High", "Low")
```

When you have a fixed number of cases and don't want to use a long chain of if-else statements, you can use `case_when()`:

```
mtcars_cyl <- mtcars |>
  mutate(cyl_category = case_when(
    cyl == 4 ~ "four",
    cyl == 6 ~ "six",
    cyl == 8 ~ "eight"
  ))
```

The `mutate()` function is used to add the new variable, and `case_when()` is used to assign the values “four”, “six”, or “eight” to the new variable based on the number of cylinders in each car. Both functions are part of the `dplyr` package (see chapter Section 7.4.1).

### 7.4.3. Examining and cleaning data with the `janitor` package

The `janitor` package follows the principles of the tidyverse and works well with the pipe operator `|>`. The `janitor` functions has many useful functions for the initial data exploration and cleaning that are essential when you load any new data set.

First, make sure the `janitor` package is installed and loaded:

## 7. Manage data

### 7.4.3.1. Clean data.frame names with `clean_names()`

I call this function frequently when I read in new data. It handles problematic variable names, especially those that are so well-preserved by `readxl::read_excel()` and `readr::read_csv()`. For example, it does the following:

- Parses letter cases and separators to a consistent format.
- Handles special characters and spaces, including transliterating characters like `æ` to `oe`.
- Appends numbers to duplicated names
- Converts “%” to “percent” and “#” to “number” to retain meaning
- Spacing (or lack thereof) around numbers is preserved

To exemplify what it does, let's create some data with awkward names and then clean them:

```
df_test <- as.data.frame(matrix(ncol = 6))
names(df_test) <- c(
  "firstName", "ábc@!*", "% successful (2009)",
  "REPEAT VALUE", "REPEAT VALUE", ""
)
df_cln <- df_test |>
  clean_names()
names(df_test)
```

```
[1] "firstName"           "ábc@!*"          "% successful (2009)"
[4] "REPEAT VALUE"        "REPEAT VALUE"      ""
```

```
names(df_cln)
```

```
[1] "first_name"          "abc"
[3] "percent_successful_2009" "repeat_value"
[5] "repeat_value_2"        "x"
```

### 7.4.3.2. Find duplicated values for specific combinations of variables with `get_dups()`

`get_dups` allows you to check for the identifying variable. In other words, it shows you duplicates for specific combinations of variables.

For example, consider the following tibble:

```
df_panel <- tibble(
  country = c(rep("a", 3), rep("b", 3), rep("c", 3)),
  year = rep(1:3, 3),
  GDP = c(1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000)
)
df_panel
```

```
# A tibble: 9 x 3
  country  year    GDP
  <chr>   <int> <dbl>
1 a         1     1000
```

## 7. Manage data

```
2 a      2 2000
3 a      3 3000
4 b      1 4000
5 b      2 5000
6 b      3 6000
7 c      1 7000
8 c      2 8000
9 c      3 9000
```

with

```
get_dupes(df_panel, country, year)
```

No duplicate combinations found of: country, year

```
# A tibble: 0 x 4
# i 4 variables: country <chr>, year <int>, dupe_count <int>, GDP <dbl>
```

we see that this is a panel dataset identified by a combination of `country` and `year`. Now let us introduce a duplicate and check again:

```
new_obs <- tibble(country = "b", year = 2, GDP = 5000)
df_panel_dup <- bind_rows(df_panel, new_obs)
get_dupes(df_panel_dup, country, year)
```

```
# A tibble: 2 x 4
  country   year dupe_count    GDP
  <chr>     <dbl>     <int> <dbl>
1 b          2        2  5000
2 b          2        2  5000
```

### 💡 Tip 11: The `plm` package

Speaking of panel datasets, it's worth mentioning the `plm` package, which is excellent for managing such data. For example, you can use `is.pbalanced` to verify whether a panel is balanced, meaning it has the same years for all countries.

```
pacman::p_load(plm)
is.pbalanced(df_panel)
```

```
[1] TRUE
```

If your panel is unbalanced, you can use `make.pbalanced` to rectify it:

```
df_unbal <- df_panel |>
  filter(row_number() != 7)
df_unbal

# A tibble: 8 x 3
```

## 7. Manage data

```
country year GDP
<chr> <int> <dbl>
1 a      1 1000
2 a      2 2000
3 a      3 3000
4 b      1 4000
5 b      2 5000
6 b      3 6000
7 c      2 8000
8 c      3 9000

is.pbalanced(df_unbal)

[1] FALSE

df_unbal_balanced <- make.pbalanced(df_unbal)
df_unbal_balanced

country year GDP
1      a   1 1000
2      a   2 2000
3      a   3 3000
4      b   1 4000
5      b   2 5000
6      b   3 6000
7      c   1   NA
8      c   2 8000
9      c   3 9000
```

### 7.4.3.3. remove\_empty() rows and columns

For cleaning Excel files that contain empty rows and columns after being read into R, `remove_empty` can be very helpful:

```
q <- data.frame(
  v1 = c(1, NA, 3),
  v2 = c(NA, NA, NA),
  v3 = c("a", NA, "b")
)
q |>
  remove_empty(c("rows", "cols"))

v1 v3
1 1 a
3 3 b
```

## 7. Manage data

### 7.4.3.4. `remove_constant()` `columns`

Removes variables from data that contain only a single constant value (with an `na.rm` option to control whether NAs should be considered as different values from the constant).

```
a <- data.frame(good = 1:3, boring = "the same")
a
```

```
good  boring
1     1 the same
2     2 the same
3     3 the same
```

```
a |>
  remove_constant()
```

```
good
1     1
2     2
3     3
```

### 7.4.4. `tabyl()` - a better version of `table()`

`tabyl()` is a tidyverse-oriented replacement for `table()`. It counts combinations of one, two, or three variables, and then can be formatted with a suite of `adorn_*` functions to look just how you want. For example:

```
mtcars |>
  tabyl(gear, cyl) |>
  adorn_totals("col") |>
  adorn_percentages("row") |>
  adorn_pct_formatting(digits = 2) |>
  adorn_ns() |>
  adorn_title()
```

gear	cyl			Total	(15)
	4	6	8		
3	6.67% (1)	13.33% (2)	80.00% (12)	100.00	
4	66.67% (8)	33.33% (4)	0.00% (0)	100.00	(12)
5	40.00% (2)	20.00% (1)	40.00% (2)	100.00% (5)	

Learn more in the [tabyls vignette](#).

## 7.5. User-defined functions and conflicts

One of the great strengths of R is the user's ability to add functions. Sometimes there is a small task (or series of tasks) you need done and you find yourself having to repeat it multiple times. In these types of situations it can be helpful to create your own custom function. The structure of a function is given below:

```
name_of_function <- function(argument1, argument2) {
  statements or code that does something
  return(something)
}
```

First you give your function a name. Then you assign value to it, where the value is the function. When defining the function you will want to provide the list of arguments required (inputs and/or options to modify behavior of the function), and wrapped between curly brackets place the tasks that are being executed on/using those arguments. The argument(s) can be any type of object (like a scalar, a matrix, a dataframe, a vector, a logical, etc), and it's not necessary to define what it is in any way. Finally, you can return the value of the object from the function, meaning pass the value of it into the global environment. The important idea behind functions is that objects that are created within the function are local to the environment of the function – they don't exist outside of the function. Note, a function doesn't require any arguments.

Let's try creating a simple example function. This function will take in a numeric value as input, and return the squared value.

```
square_it <- function(x) {
  square <- x * x
  return(square)
}
```

Now, we can use the function as we would any other function. We type out the name of the function, and inside the parentheses we provide a numeric value `x`:

```
square_it(5)
```

```
[1] 25
```

Let us get back to script with sales and try to calculate the monthly growth rates of revenue using a self-written function.

The formula of a growth rate is clear:

$$g = \left( \frac{y_t - y_{t-1}}{y_{t-1}} \right) \cdot 100 = \left( \frac{y_t}{y_{t-1}} - 1 \right) \cdot 100$$

So the challenge is to divide the value of `revenue` with the value of the previous period, a.k.a. the lagged value. Let us assume that the function `lag()` can give you exactly that value of a vector. Lets try it out:

## 7. Manage data

```
lag(revenue)
```

```
[1] 0 700 1400 350 175 28 56 0 0 0 0 0  
attr(,"tsp")  
[1] 0 11 1
```

```
(revenue/lag(revenue)-1)*100
```

```
[1] NaN 0 0 0 0 0 NaN NaN NaN NaN NaN  
attr(,"tsp")  
[1] 0 11 1
```

Unfortunately, this does not work out. The `lag()` function does not work as we think it should. Well, the reason is simply that we are using the wrong function. The current `lag()` function is part of the `stats` package which is part of the package `stats` which is part of R base and is loaded automatically. The `lag()` function we aim to use stems from the `dplyr` package which we must install and load to be able to use it. So let's do it:

```
# check if the package is installed  
find.package("dplyr")
```

```
[1] "/usr/local/lib/R/site-library/dplyr"
```

```
# I already installed the package so I can just load it  
# install.packages("dplyr")  
library("dplyr")
```

```
Attaching package: 'dplyr'
```

```
The following objects are masked from 'package:plm':
```

```
between, lag, lead
```

```
The following objects are masked from 'package:stats':
```

```
filter, lag
```

```
The following objects are masked from 'package:base':
```

```
intersect, setdiff, setequal, union
```

This message informs us that among other functions the `lag()` function is *masked*. That means that now the function of the newly loaded package is active. This is one example of why I highly recommend to unload all packages at the beginning of a script and then to use `p_load` to install and load the packages that should be used in the upcoming script:

## 7. Manage data

```
pacman::p_unload(all)
pacman::p_load(dplyr)
```

So, let's try again:

```
lag(revenue)

[1] NA 0 700 1400 350 175 28 56 0 0 0 0

(revenue/lag(revenue)-1)*100

[1] NA Inf 100 -75 -50 -84 100 -100 NaN NaN NaN NaN
```

That looks good now. And here is a way to calculate growth rates with a self-written function:

```
growth_rate <- function(x) {
  (x / lag(x) - 1) * 100
}
growth_rate(revenue)

[1] NA Inf 100 -75 -50 -84 100 -100 NaN NaN NaN NaN

sales_gr_rate <- growth_rate(revenue)
sales_gr_rate

[1] NA Inf 100 -75 -50 -84 100 -100 NaN NaN NaN NaN
```

In R, all functions are written by users, and it is not uncommon for two people to name their functions identically. In such cases, we must resolve the conflict by choosing which function to use. To use the lag function from the `stats` package, you can use the double colon operator `::` like this `stats::lag()`.

## 7.6. Example: How to explore a dataset

```
# Creating dataframe
df <- tibble(
  integer_var, numeric_var, character_var, factor_var, logical_var, date_var,
)

# Overview of the data
head(df)
```

## 7. Manage data

```
# A tibble: 5 x 6
  integer_var numeric_var character_var factor_var logical_var date_var
    <dbl>      <dbl> <chr>        <fct>     <lgl>      <date>
1       1        1.1 apple       red       TRUE 2022-01-01
2       2        2.2 banana     yellow    TRUE 2022-02-01
3       3        NA orange     red       TRUE 2022-03-01
4       4        4.4 cherry    blue      FALSE 2022-04-01
5       5        5.5 grape     green     TRUE 2022-05-01
```

```
summary(df)
```

```
integer_var  numeric_var   character_var   factor_var logical_var
Min.    :1   Min.    :1.100  Length:5       blue    :1   Mode :logical
1st Qu.:2   1st Qu.:1.925  Class :character  green   :1   FALSE:1
Median :3   Median :3.300  Mode   :character  red    :2   TRUE  :4
Mean   :3   Mean    :3.300
3rd Qu.:4   3rd Qu.:4.675
Max.   :5   Max.    :5.500
NA's    :1

date_var
Min.    :2022-01-01
1st Qu.:2022-02-01
Median :2022-03-01
Mean   :2022-03-02
3rd Qu.:2022-04-01
Max.   :2022-05-01
```

```
glimpse(df)
```

```
Rows: 5
Columns: 6
$ integer_var  <dbl> 1, 2, 3, 4, 5
$ numeric_var   <dbl> 1.1, 2.2, NA, 4.4, 5.5
$ character_var <chr> "apple", "banana", "orange", "cherry", "grape"
$ factor_var    <fct> red, yellow, red, blue, green
$ logical_var   <lgl> TRUE, TRUE, TRUE, FALSE, TRUE
$ date_var      <date> 2022-01-01, 2022-02-01, 2022-03-01, 2022-04-01, 2022-05-~
```

```
# look closer at variables
```

```
# unique values
```

```
unique(df$integer_var)
```

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

```
unique(df$factor_var)
```

## 7. Manage data

```
[1] red     yellow blue   green  
Levels: blue green red yellow
```

```
table(df$factor_var)
```

	blue	green	red	yellow
1	1	2	1	

```
length(unique(df$factor_var))
```

```
[1] 4
```

```
# distributions  
df |> count(factor_var)
```

```
# A tibble: 4 x 2  
  factor_var     n  
  <fct>       <int>  
1 blue          1  
2 green         1  
3 red           2  
4 yellow        1
```

```
prop.table(table(df$factor_var))
```

	blue	green	red	yellow
	0.2	0.2	0.4	0.2

```
df |>  
  count(factor_var) |>  
  mutate(prop = n / sum(n))
```

```
# A tibble: 4 x 3  
  factor_var     n   prop  
  <fct>       <int> <dbl>  
1 blue          1   0.2  
2 green         1   0.2  
3 red           2   0.4  
4 yellow        1   0.2
```

```
aggregate(df$numeric_var,  
  by = list(fruit = df$factor_var),  
  mean  
)
```

## 7. Manage data

```
fruit   x
1  blue 4.4
2 green 5.5
3   red NA
4 yellow 2.2
```

```
# --> the mean of red cannot be calculated as there is a NA in it
# Solution: exclude NAs from calculation:
aggregate(df$numeric_var,
  by = list(fruit = df$factor_var),
  mean,
  na.rm = TRUE
)
```

```
fruit   x
1  blue 4.4
2 green 5.5
3   red 1.1
4 yellow 2.2
```

```
# install.packages("janitor")
require("janitor")
mtcars |>
  tabyl(cyl)
```

```
cyl  n percent
  4 11 0.34375
  6  7 0.21875
  8 14 0.43750
```

```
mtcars |>
  tabyl(cyl, hp)
```

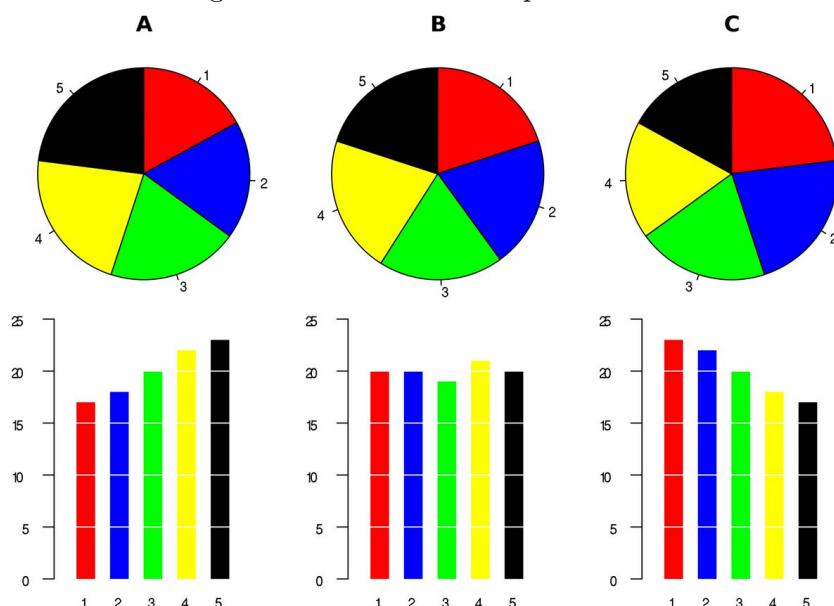
```
cyl 52 62 65 66 91 93 95 97 105 109 110 113 123 150 175 180 205 215 230 245
  4  1  1  1  2  1  1  1  1  0  1  0  1  0  0  0  0  0  0  0  0  0  0  0
  6  0  0  0  0  0  0  0  0  1  0  3  0  2  0  1  0  0  0  0  0  0  0
  8  0  0  0  0  0  0  0  0  0  0  0  0  0  2  2  3  1  1  1  1  2
264 335
  0  0
  0  0
  1  1
```

## 8. Visualize data

Data visualization is an art. The purposes of visualizing data are manifold. You can emphasize facts, get known to data, detect anomalies, and communicate a large amount of information simply and intuitive. Whatever your goal is, thousand of appropriate ways exist to visualize data. Many decisions to take are simply a matter of taste. However, there are some conventions and guidelines that help you to make on average better decisions when designing a visualization:

- Good graphs are easy to understand and eye catching.
- Graphs can be misleading and manipulative and that is opposing to the ideas of science. Thus, be responsible and honest.
- Minimize colors and other attention-grabbing elements that are not directly related to the data of interest. Worldwide, there are approximately 300 million color blind people. In particular, red, green or blue light are problematic to color blind people. Thus, better rely on color schemes that are designed for colorblind people.
- Don't truncate an axis or change the scaling within an axis just to make you your story more appealing. Show the full scale of the graph, then zoom to show the data of interest, if necessary.
- Label and describe your chart sufficiently so that everybody can fully understand the content of the shown data set and statistics without having to study the notes of the graph for too long.
- Don't do pie charts. They may look simple, but they're tricky to get right and there are usually better alternatives. Humans are not very good at comparing the size of angles and as there's no scale in pie plots, reading accurate values is difficult. Figure Figure 8.1 may proof this.

Figure 8.1.: Pie charts are problematic



Source: [https://en.wikipedia.org/wiki/Pie\\_chart](https://en.wikipedia.org/wiki/Pie_chart)

 More tips

- [Data Visualization: Chart Dos and Don'ts \(by Duke University\)](#)
- [Graphs and Visualising Data](#) by Oliver Kirchkamp. In particular, I highly recommend his [handout \[Kirchkamp, 2018\]](#). It discusses many pitfalls of visualizing data, instructs how to do good graphs, and he shows the corresponding R code of all graphs.
- The [From Data to Viz](#) website leads you to the most appropriate graph for your data. It links to the code to build it and lists common caveats you should avoid.
- The [R Graph Gallery](#) and [R CHARTS by R CODER](#) shows graphs and the corresponding R code to replicate the graphs
- The work of [Edward Tufte](#) and his book *The Visual Display of Quantitative Information* [[Tufte, 2022](#)] are classical readings.

A great resource to learn how to visualize data is [Wickham and Grolemund \[2023\]](#). As I cannot do that any better, I refer to that source and refrain from writing section myself. It introduces the `ggplot` function which is part of the `ggplot2` package which, in turn, is part of the tidyverse package. Thus, if you've installed and loaded tidyverse, you automatically have access to `ggplot`. Creating beautiful and informative graphs is easy with `ggplot`. To proof that claim, study the chapter ([Data visualization](#)) of [Wickham and Grolemund \[2023\]](#). Another good resource on modern data visualization is [Kabacoff \[2024\]](#).

# 9. Collection of exercises

## 9.1. Import data with `c()`

Table 9.1 shows COVID for three states in Germany:

Table 9.1.: Covid cases and deaths till August 2022

state	Bavaria	North Rhine-Westphalia	Baden-Württemberg
deaths (in mio)	4,92M	5,32M	3,69M
cases	24.111	25.466	16.145

Write down the code you would need to put into the R-console...

- ...to store each of variables `state` and `deaths` in a vector.
- ...to store both vectors in a data frame with the name `df_covid`.
- ...to store both vectors in a tibble with the name `tbl_covid`.

### Solution

The script uses the following functions: `c`, `data.frame`, `tibble`.

#### R script

```
# Solution to excercise "Import data":\n\n# load packages\nif (!require(pacman)) install.packages("pacman")\npacman::p_load(tibble)\n\nstate <- c("BY", "NRW", "BW")\ndeaths <- c(4.92, 5.32, 3.69)\ncases <- c(24111, 25466, 16145)\ndf_covid <- data.frame(state, deaths)\ntbl_covid <- tibble(state, deaths)\n\ntbl_covid\n\nsuppressMessages(pacman::p_unload(tibble))
```

#### Output of the R script

```
# Solution to excercise "Import data":

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tibble)

state <- c("BY", "NRW", "BW")
deaths <- c(4.92, 5.32, 3.69)
cases <- c(24111, 25466, 16145)
df_covid <- data.frame(state, deaths)
tbl_covid <- tibble(state, deaths)

tbl_covid

# A tibble: 3 x 2
  state   deaths
  <chr>   <dbl>
1 BY      4.92
2 NRW     5.32
3 BW      3.69

suppressMessages(pacman::p_unload(tibble))
```

## 9.2. Filter and select observations

Set up R, RStudio, and R packages

Open [this interactive tutorial](#) and work through it.

The script uses among others the following functions: `filter`, `is.na`, `select`.

## 9.3. Base R,`%in%` operator, and the pipe `|>`

- Using the `mtcars` dataset, write code to create a new dataframe that includes only the rows where the number of cylinders is either 4 or 6, and the weight (`wt`) is less than 3.5.

Do this in two different ways using:

1. The `%in%` operator and the pipe `|>`.
2. Base R without the pipe `|>`.

Compare the resulting dataframes using the `identical()` function.

- b) Using the `mtcars` dataset, generate a logical variable that indicates with TRUE all cars with either 4 or 6 cylinders that `wt` is less than 3.5 and add this variable to a new dataset.

 Solution

The script uses the following functions: `c`, `filter`, `identical`, `if_else`, `mutate`, `subset`, `transform`, `with`.

**R script**

## 9. Collection of exercises

```
# Base R or pipe
# exe_base_pipe.R
# Stephan Huber; 2023-05-08

# setwd("/home/sthu/Dropbox/hsf/test")
rm(list=ls())

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(datasets, tidyverse)

# a)
# Using the pipe |>
# Select rows where cyl is 4 or 6 and wt is less than 3.5
df1 <- mtcars |>
  filter(cyl %in% c(4, 6) & wt < 3.5)
df1

# Without the pipe |>
# Select rows where cyl is 4 or 6 and wt is less than 3.5
df2 <- subset(mtcars, cyl %in% c(4, 6) & wt < 3.5)
df2

# Check if the resulting dataframe is identical to the expected output
identical(df1, df2)

# b)
# Using the pipe |> and tidyverse (mutate)
df3 <- mtcars |>
  mutate(cyl_4_or_6 =
    if_else(cyl %in% c(4, 6) & wt < 3.5, TRUE, FALSE))
df3

# without pipe and with base R (transform)
df4 <- mtcars
df4$cyl_4_or_6 <- with(mtcars, cyl %in% c(4, 6) & wt < 3.5)

# Alternatively in one line:
df5 <- transform(mtcars, cyl_4_or_6 = cyl %in% c(4,6) & wt < 3.5)

# Check if the resulting dataframe is identical to the expected output
identical(df3, df4)
identical(df3, df5)

# unload packages
suppressMessages(pacman::p_unload(datasets, tidyverse))
```

### Output of the R script

## 9. Collection of exercises

```
# Base R or pipe
# exe_base_pipe.R
# Stephan Huber; 2023-05-08

# setwd("/home/sthu/Dropbox/hsf/test")
rm(list=ls())

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(datasets, tidyverse)

# a)
# Using the pipe |>
# Select rows where cyl is 4 or 6 and wt is less than 3.5
df1 <- mtcars |>
  filter(cyl %in% c(4, 6) & wt < 3.5)
df1
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1
Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3	1
Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4	4
Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	0	4	4
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	1	4	1
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4	2
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	1	4	1
Toyota Corona	21.5	4	120.1	97	3.70	2.465	20.01	1	0	3	1
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	1	4	1
Porsche 914-2	26.0	4	120.3	91	4.43	2.140	16.70	0	1	5	2
Lotus Europa	30.4	4	95.1	113	3.77	1.513	16.90	1	1	5	2
Ferrari Dino	19.7	6	145.0	175	3.62	2.770	15.50	0	1	5	6
Volvo 142E	21.4	4	121.0	109	4.11	2.780	18.60	1	1	4	2

```
# Without the pipe |>
# Select rows where cyl is 4 or 6 and wt is less than 3.5
df2 <- subset(mtcars, cyl %in% c(4, 6) & wt < 3.5)
df2
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1
Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3	1

## 9. Collection of exercises

Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4	4
Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	0	4	4
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	1	4	1
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4	2
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	1	4	1
Toyota Corona	21.5	4	120.1	97	3.70	2.465	20.01	1	0	3	1
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	1	4	1
Porsche 914-2	26.0	4	120.3	91	4.43	2.140	16.70	0	1	5	2
Lotus Europa	30.4	4	95.1	113	3.77	1.513	16.90	1	1	5	2
Ferrari Dino	19.7	6	145.0	175	3.62	2.770	15.50	0	1	5	6
Volvo 142E	21.4	4	121.0	109	4.11	2.780	18.60	1	1	4	2

```
# Check if the resulting dataframe is identical to the expected output
identical(df1, df2)
```

[1] TRUE

```
# b)
# Using the pipe |> and tidyverse (mutate)
df3 <- mtcars |>
  mutate(cyl_4_or_6 =
         if_else(cyl %in% c(4, 6) & wt < 3.5, TRUE, FALSE))
df3
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1
Hornet Sportabout	18.7	8	360.0	175	3.15	3.440	17.02	0	0	3	2
Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3	1
Duster 360	14.3	8	360.0	245	3.21	3.570	15.84	0	0	3	4
Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4	4
Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	0	4	4
Merc 450SE	16.4	8	275.8	180	3.07	4.070	17.40	0	0	3	3
Merc 450SL	17.3	8	275.8	180	3.07	3.730	17.60	0	0	3	3
Merc 450SLC	15.2	8	275.8	180	3.07	3.780	18.00	0	0	3	3
Cadillac Fleetwood	10.4	8	472.0	205	2.93	5.250	17.98	0	0	3	4
Lincoln Continental	10.4	8	460.0	215	3.00	5.424	17.82	0	0	3	4
Chrysler Imperial	14.7	8	440.0	230	3.23	5.345	17.42	0	0	3	4
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	1	4	1
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4	2
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	1	4	1
Toyota Corona	21.5	4	120.1	97	3.70	2.465	20.01	1	0	3	1
Dodge Challenger	15.5	8	318.0	150	2.76	3.520	16.87	0	0	3	2

9. Collection of exercises

AMC Javelin	15.2	8	304.0	150	3.15	3.435	17.30	0	0	3	2
Camaro Z28	13.3	8	350.0	245	3.73	3.840	15.41	0	0	3	4
Pontiac Firebird	19.2	8	400.0	175	3.08	3.845	17.05	0	0	3	2
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	1	4	1
Porsche 914-2	26.0	4	120.3	91	4.43	2.140	16.70	0	1	5	2
Lotus Europa	30.4	4	95.1	113	3.77	1.513	16.90	1	1	5	2
Ford Pantera L	15.8	8	351.0	264	4.22	3.170	14.50	0	1	5	4
Ferrari Dino	19.7	6	145.0	175	3.62	2.770	15.50	0	1	5	6
Maserati Bora	15.0	8	301.0	335	3.54	3.570	14.60	0	1	5	8
Volvo 142E	21.4	4	121.0	109	4.11	2.780	18.60	1	1	4	2
			cyl_4_or_6								
Mazda RX4			TRUE								
Mazda RX4 Wag			TRUE								
Datsun 710			TRUE								
Hornet 4 Drive			TRUE								
Hornet Sportabout			FALSE								
Valiant			TRUE								
Duster 360			FALSE								
Merc 240D			TRUE								
Merc 230			TRUE								
Merc 280			TRUE								
Merc 280C			TRUE								
Merc 450SE			FALSE								
Merc 450SL			FALSE								
Merc 450SLC			FALSE								
Cadillac Fleetwood			FALSE								
Lincoln Continental			FALSE								
Chrysler Imperial			FALSE								
Fiat 128			TRUE								
Honda Civic			TRUE								
Toyota Corolla			TRUE								
Toyota Corona			TRUE								
Dodge Challenger			FALSE								
AMC Javelin			FALSE								
Camaro Z28			FALSE								
Pontiac Firebird			FALSE								
Fiat X1-9			TRUE								
Porsche 914-2			TRUE								
Lotus Europa			TRUE								
Ford Pantera L			FALSE								
Ferrari Dino			TRUE								
Maserati Bora			FALSE								
Volvo 142E			TRUE								

```
# without pipe and with base R (transform)
df4 <- mtcars
df4$cyl_4_or_6 <- with(mtcars, cyl %in% c(4, 6) & wt < 3.5)

# Alternatively in one line:
df5 <- transform(mtcars, cyl_4_or_6 = cyl %in% c(4,6) & wt < 3.5)

# Check if the resulting dataframe is identical to the expected output
identical(df3, df4)
```

[1] TRUE

```
identical(df3, df5)
```

[1] TRUE

```
# unload packages
suppressMessages(pacman::p_unload(datasets, tidyverse))
```

## 9.4. Generate and drop variables

Use the *mtcars* dataset. It is part of the package *datasets* and can be called with

```
mtcars
```

- Create a new tibble called `mtcars_new` using the pipe operator `|>`. Generate a new dummy variable called `d_cyl_6to8` that takes the value 1 if the number of cylinders (`cyl`) is greater than 6, and 0 otherwise. Do all of this in a single pipe.
- Generate a new dummy variable called `posercar` that takes a value of 1 if a car has more than 6 cylinders (`cyl`) and can drive less than 18 miles per gallon (`mpg`), and 0 otherwise. Add this variable to the tibble `mtcars_new`.
- Remove the variable `d_cyl_6to8` from the data frame.

### Solution

The script uses the following functions: `as_tibble`, `if_else`, `mutate`, `rownames_to_column`, `select`.

### R script

```

# Generate and drop variables
# exe_genanddrop.R
# Stephan Huber; 2023-05-09

# setwd("/home/sthu/Dropbox/hsf/test")
rm(list = ls())

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(datasets, tidyverse)

# a)
mtcars_new <- mtcars |>
  rownames_to_column(var = "car") |>
  as_tibble() |>
  mutate(d_cyl_6to8 = if_else(cyl > 6, 1, 0))
mtcars_new

# b)
mtcars_new <- mtcars_new |>
  mutate(posercar = if_else(cyl > 6 & mpg < 18, 1, 0))
mtcars_new

# c)
mtcars_new <- mtcars_new |>
  select(-d_cyl_6to8)
mtcars_new

# unload packages
suppressMessages(pacman::p_unload(datasets, tidyverse))

```

### Output of the R script

```

# Generate and drop variables
# exe_genanddrop.R
# Stephan Huber; 2023-05-09

# setwd("/home/sthu/Dropbox/hsf/test")
rm(list = ls())

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(datasets, tidyverse)

# a)
mtcars_new <- mtcars |>
  rownames_to_column(var = "car") |>
  as_tibble() |>
  mutate(d_cyl_6to8 = if_else(cyl > 6, 1, 0))
mtcars_new

```

## 9. Collection of exercises

```
# A tibble: 32 x 13
  car      mpg cyl disp   hp drat    wt  qsec   vs   am gear carb
  <chr>     <dbl> <dbl>
1 Mazda RX4  21     6   160   110  3.9   2.62  16.5   0     1   4    4
2 Mazda RX4~ 21     6   160   110  3.9   2.88  17.0   0     1   4    4
3 Datsun 710 22.8   4   108   93   3.85  2.32  18.6   1     1   4    1
4 Hornet 4 D~ 21.4   6   258   110  3.08  3.22  19.4   1     0   3    1
5 Hornet Spo~ 18.7   8   360   175  3.15  3.44  17.0   0     0   3    2
6 Valiant     18.1   6   225   105  2.76  3.46  20.2   1     0   3    1
7 Duster 360 14.3   8   360   245  3.21  3.57  15.8   0     0   3    4
8 Merc 240D   24.4   4   147.   62   3.69  3.19  20.0   1     0   4    2
9 Merc 230    22.8   4   141.   95   3.92  3.15  22.9   1     0   4    2
10 Merc 280   19.2   6   168.  123   3.92  3.44  18.3   1    0   4    4
# i 22 more rows
# i 1 more variable: d_cyl_6to8 <dbl>
```

```
# b)
mtcars_new <- mtcars_new |>
  mutate(posercar = if_else(cyl > 6 & mpg < 18, 1, 0))
mtcars_new
```

```
# A tibble: 32 x 14
  car      mpg cyl disp   hp drat    wt  qsec   vs   am gear carb posercar
  <chr>     <dbl> <dbl>
1 Mazda RX4  21     6   160   110  3.9   2.62  16.5   0     1   4    4     0
2 Mazda RX4~ 21     6   160   110  3.9   2.88  17.0   0     1   4    4     0
3 Datsun 710 22.8   4   108   93   3.85  2.32  18.6   1     1   4    1     0
4 Hornet 4 D~ 21.4   6   258   110  3.08  3.22  19.4   1     0   3    1     0
5 Hornet Spo~ 18.7   8   360   175  3.15  3.44  17.0   0     0   3    2     0
6 Valiant     18.1   6   225   105  2.76  3.46  20.2   1     0   3    1     0
7 Duster 360 14.3   8   360   245  3.21  3.57  15.8   0     0   3    4     0
8 Merc 240D   24.4   4   147.   62   3.69  3.19  20.0   1     0   4    2     0
9 Merc 230    22.8   4   141.   95   3.92  3.15  22.9   1     0   4    2     0
10 Merc 280   19.2   6   168.  123   3.92  3.44  18.3   1    0   4    4     0
# i 22 more rows
# i 2 more variables: d_cyl_6to8 <dbl>, posercar <dbl>
```

```
# c)
mtcars_new <- mtcars_new |>
  select(-d_cyl_6to8)
mtcars_new
```

```
# A tibble: 32 x 13
  car      mpg cyl disp   hp drat    wt  qsec   vs   am gear carb
  <chr>     <dbl> <dbl>
1 Mazda RX4  21     6   160   110  3.9   2.62  16.5   0     1   4    4
2 Mazda RX4~ 21     6   160   110  3.9   2.88  17.0   0     1   4    4
3 Datsun 710 22.8   4   108   93   3.85  2.32  18.6   1     1   4    1
4 Hornet 4 D~ 21.4   6   258   110  3.08  3.22  19.4   1     0   3    1
```

```

      5 Hornet Spo~ 18.7     8   360     175  3.15  3.44  17.0     0     0    3    2
      6 Valiant     18.1     6   225     105  2.76  3.46  20.2     1     0    3    1
      7 Duster 360 14.3     8   360     245  3.21  3.57  15.8     0     0    3    4
      8 Merc 240D   24.4     4   147.     62   3.69  3.19  20       1     0    4    2
      9 Merc 230    22.8     4   141.     95   3.92  3.15  22.9     1     0    4    2
     10 Merc 280   19.2     6   168.    123   3.92  3.44  18.3     1     0    4    4
# i 22 more rows
# i 1 more variable: posercar <dbl>

# unload packages
suppressMessages(pacman::p_unload(datasets, tidyverse))

```

## 9.5. Subsetting

1. Check to see if you have the `mtcars` dataset by entering `mtcars`.
2. Save the `mtcars` dataset in an object named `cars`.
3. What class is `cars`?
4. How many observations (rows) and variables (columns) are in the `mtcars` dataset?
5. Rename `mpg` in `cars` to `MPG`. Use `rename()`.
6. Convert the column names of `cars` to all upper case. Use `rename_all`, and the `toupper` function.
7. Convert the rownames of `cars` to a column called `car` using `rownames_to_column`.
8. Subset the columns from `cars` that end in “p” and call it `pvars` using `ends_with()`.
9. Create a subset `cars` that only contains the columns: `wt`, `qsec`, and `hp` and assign this object to `carsSub`. (Use `select()`.)
10. What are the dimensions of `carsSub`? (Use `dim()`.)
11. Convert the column names of `carsSub` to all upper case. Use `rename_all()`, and `toupper()` (or `colnames()`).
12. Subset the rows of `cars` that get more than 20 miles per gallon (`mpg`) of fuel efficiency. How many are there? (Use `filter()`.)
13. Subset the rows that get less than 16 miles per gallon (`mpg`) of fuel efficiency and have more than 100 horsepower (`hp`). How many are there? (Use `filter()` and the pipe operator.)
14. Create a subset of the `cars` data that only contains the columns: `wt`, `qsec`, and `hp` for cars with 8 cylinders (`cyl`) and reassign this object to `carsSub`. What are the dimensions of this dataset? Do not use the pipe operator.
15. Create a subset of the `cars` data that only contains the columns: `wt`, `qsec`, and `hp` for cars with 8 cylinders (`cyl`) and reassign this object to `carsSub2`. Use the pipe operator.
16. Re-order the rows of `carsSub` by weight (`wt`) in increasing order. (Use `arrange()`.)
17. Create a new variable in `carsSub` called `wt2`, which is equal to  $wt^2$ , using `mutate()` and piping `|>`.

### Solution

The script uses the following functions: `arrange`, `class`, `dim`, `ends_with`, `filter`, `mutate`, `ncol`, `nrow`, `rename`, `rename_all`, `rownames_to_column`, `select`.

### R script

## 9. Collection of exercises

```
# Subsetting with \R
# exe_subset.R
# Stephan Huber; 2022-06-07

# setwd("/home/sthu/Dropbox/hsf/22-ss/dsda/work/")
rm(list = ls())

# 0
# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, dplyr, tibble)

# 1
mtcars

# 2
cars <- mtcars

# 3
class(cars)

# 4
dim(cars)

# Alternative
ncol(cars)
nrow(cars)

# 5
cars <- rename(cars, MPG = mpg)

# 6
cars <- rename_all(cars, toupper)
# if you like lower cases:
# cars <- rename_all(cars, tolower)

# 7
cars <- rownames_to_column(mtcars, var = "car")

# 8
pvars <- select(cars, car, ends_with("p"))

# 9
carsSub <- select(cars, car, wt, qsec, hp)

# 10
dim(carsSub)

# 11
carsSub <- rename_all(carsSub, toupper)

# 12
cars_mpg <- filter(cars, mpg > 20)
dim(cars_mpg)
```

## Output of the R script

```
# Subsetting with \R
# exe_subset.R
# Stephan Huber; 2022-06-07

# setwd("/home/sthu/Dropbox/hsf/22-ss/dsda/work/")
rm(list = ls())

# 0
# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, dplyr, tibble)

# 1
mtcars
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1
Hornet Sportabout	18.7	8	360.0	175	3.15	3.440	17.02	0	0	3	2
Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3	1
Duster 360	14.3	8	360.0	245	3.21	3.570	15.84	0	0	3	4
Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4	4
Merc 280C	17.8	6	167.6	123	3.92	3.440	18.90	1	0	4	4
Merc 450SE	16.4	8	275.8	180	3.07	4.070	17.40	0	0	3	3
Merc 450SL	17.3	8	275.8	180	3.07	3.730	17.60	0	0	3	3
Merc 450SLC	15.2	8	275.8	180	3.07	3.780	18.00	0	0	3	3
Cadillac Fleetwood	10.4	8	472.0	205	2.93	5.250	17.98	0	0	3	4
Lincoln Continental	10.4	8	460.0	215	3.00	5.424	17.82	0	0	3	4
Chrysler Imperial	14.7	8	440.0	230	3.23	5.345	17.42	0	0	3	4
Fiat 128	32.4	4	78.7	66	4.08	2.200	19.47	1	1	4	1
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4	2
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.90	1	1	4	1
Toyota Corona	21.5	4	120.1	97	3.70	2.465	20.01	1	0	3	1
Dodge Challenger	15.5	8	318.0	150	2.76	3.520	16.87	0	0	3	2
AMC Javelin	15.2	8	304.0	150	3.15	3.435	17.30	0	0	3	2
Camaro Z28	13.3	8	350.0	245	3.73	3.840	15.41	0	0	3	4
Pontiac Firebird	19.2	8	400.0	175	3.08	3.845	17.05	0	0	3	2
Fiat X1-9	27.3	4	79.0	66	4.08	1.935	18.90	1	1	4	1
Porsche 914-2	26.0	4	120.3	91	4.43	2.140	16.70	0	1	5	2
Lotus Europa	30.4	4	95.1	113	3.77	1.513	16.90	1	1	5	2
Ford Pantera L	15.8	8	351.0	264	4.22	3.170	14.50	0	1	5	4
Ferrari Dino	19.7	6	145.0	175	3.62	2.770	15.50	0	1	5	6
Maserati Bora	15.0	8	301.0	335	3.54	3.570	14.60	0	1	5	8
Volvo 142E	21.4	4	121.0	109	4.11	2.780	18.60	1	1	4	2

## 9. Collection of exercises

```
# 2
cars <- mtcars

# 3
class(cars)

[1] "data.frame"

# 4
dim(cars)

[1] 32 11

# Alternative
ncol(cars)

[1] 11

nrow(cars)

[1] 32

# 5
cars <- rename(cars, MPG = mpg)

# 6
cars <- rename_all(cars, toupper)
# if you like lower cases:
# cars <- rename_all(cars, tolower)

# 7
cars <- rownames_to_column(mtcars, var = "car")

# 8
pvars <- select(cars, car, ends_with("p"))

# 9
carsSub <- select(cars, car, wt, qsec, hp)

# 10
dim(carsSub)

[1] 32 4

# 11
carsSub <- rename_all(carsSub, toupper)

# 12
cars_mpg <- filter(cars, mpg > 20)
dim(cars_mpg)
```

```
[1] 14 12
```

```
# 13
cars_whatever <- filter(cars, mpg < 16 & hp > 100)

# 14
carsSub <- filter(cars, cyl == 8)
carsSub <- select(carsSub, wt, qsec, hp, car)
dim(carsSub)
```

```
[1] 14 4
```

```
# 15
# Alternative with pipe operator:
carsSub <- cars |>
  filter(cyl == 8) |>
  select(wt, qsec, hp, car)

# 16
carsSub <- arrange(carsSub, wt)

# 17
carsSub <- carsSub |>
  mutate(wt2 = wt^2)

# Alternatively you can put everything into one pipe:
carsSub2 <- cars |>
  filter(cyl == 8) |>
  select(wt, qsec, hp, car) |>
  arrange(carsSub, wt) |>
  mutate(wt2 = wt^2)

# unload packages
suppressMessages(pacman::p_unload(tidyverse, dplyr, tibble))
```

## 9.6. Weighting, data management and regression

Consider the data of Table Table 9.2 and solve the following exercises:

- a) Add variable `misone` that is 1 if there is a missing and 0 otherwise. (*Hint: Use `case_when` and `is.na()`.*)

```
library(dplyr)
```

```
df <- df |>
```

Table 9.2.: Data

	v1	v2	v3	v4
	1	A	NA	
2	NA	0.2	0.4	
	C	0.3	0.1	
4	D	NA	3	
5	E	0.5	1.5	

	v1	v2	v3	v4	misone
	1	A	NA		1
2	NA	0.2	0.4	1	
	C	0.3	0.1	1	
4	D	NA	3	1	
5	E	0.5	1.5	0	

```

mutate(misone = case_when(
  v1 == "" | is.na(v1) ~ 1,
  v2 == "" | is.na(v2) ~ 1,
  v3 == "" | is.na(v3) ~ 1,
  v4 == "" | is.na(v4) ~ 1,
  TRUE ~ 0
))

tt(df, output="markdown")

```

- b) Add variable `miscount` that counts how many observations are missing in each row. (*Hint: Use `mutate_all`, `rowSums`, and `pick(everything())`*)

```

test_df <- df |>
  mutate_all(~ if_else(is.na(.) | . == "", 1, 0)) |>
  mutate(miscount = rowSums(pick(everything())))

test_df_miscount <- test_df |>
  select(miscount)

df <- bind_cols(df, test_df_miscount)

tt(df, output="markdown")

```

- c) Use the function `rowwise` to calculate the `NA` and `" "` observations. (*Hint: Use `is.na` and `pick(everything())`.*)

## 9. Collection of exercises

v1	v2	v3	v4	misone	miscount
1	A	NA		1	2
2	NA	0.2	0.4	1	1
	C	0.3	0.1	1	1
4	D	NA	3	1	1
5	E	0.5	1.5	0	0

v1	v2	v3	v4	misone	miscount	countNA	countOK
1	A	NA		1	2	1	1
2	NA	0.2	0.4	1	1	1	0
	C	0.3	0.1	1	1	0	1
4	D	NA	3	1	1	1	0
5	E	0.5	1.5	0	0	0	0

```
df <- df |>
  rowwise() |>
  mutate(countNA = sum(is.na(pick(everything())))) |>
  mutate(countOK = sum(pick(everything()) == "", na.rm = TRUE)) |>
  ungroup()

tt(df, output="markdown")
```

- d) Add variable `mispercent` that measures the percentage of missings and a variable `mis30up` that is 1 if the percentage is above 30%. (*Hint: Use `mutate`, `select`, `ifelse`, and `bind_cols`.*)

```
test_df_mis30up <- test_df |>
  mutate(fraction = miscount / 4) |>
  mutate(mis30up = ifelse(fraction > 0.3, 1, 0)) |>
  select(mis30up, fraction)

df <- bind_cols(df, test_df_mis30up)

tt(df, output="markdown")
```

- e) Calculate the average of the numeric variables `v1`, `v3`, and `v4`. Name the variable `average`. (*Hint: Use `as.numeric`, `rowwise`, and `mean`.*)

v1	v2	v3	v4	misone	miscount	countNA	countOK	mis30up	fraction
1	A	NA		1	2	1	1	1	0.50
2	NA	0.2	0.4	1	1	1	0	0	0.25
	C	0.3	0.1	1	1	0	1	0	0.25
4	D	NA	3	1	1	1	0	0	0.25
5	E	0.5	1.5	0	0	0	0	0	0.00

v1	v2	v3	v4	misone	miscount	countNA	countOK	miss30up	fraction	average
1	A	NA	NA	1	2	1	1	1	0.50	1.0000000
2	NA	0.2	0.4	1	1	1	0	0	0.25	0.8666667
NA	C	0.3	0.1	1	1	0	1	0	0.25	0.2000000
4	D	NA	3.0	1	1	1	0	0	0.25	3.5000000
5	E	0.5	1.5	0	0	0	0	0	0.00	2.3333333

```
df <- df |>
  mutate(
    v1 = as.numeric(v1),
    v4 = as.numeric(v4)
  )
df <- df |>
  rowwise() |>
  mutate(average = mean(c(v1, v3, v4), na.rm = TRUE)) |>
  ungroup()

test_df <- df |>
  select(v1, v3, v4, average)

tt(df, output="markdown")
```

## 9.7. Consumer prices over time

1. Read in the following data:

<https://raw.githubusercontent.com/hubchev/courses/main/dta/PCEPI.csv>

The data stem from the [Federal Reserve Bank of St. Louis \(FRED\)](#).



Tip

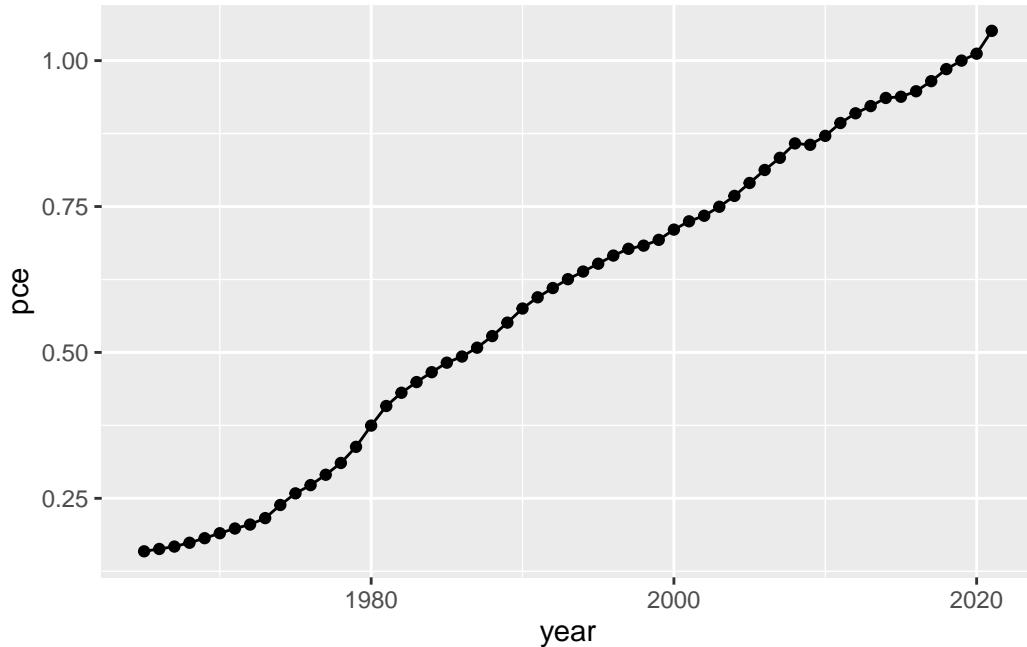
You should always try to use most recent data and you should not work outside of R. Thus, it would be optimal to download the data directly from FRED and using a R package. If this would be serious research, I would not recommend to use the data that I have downloaded from FRED, I'd recommend to use the [fredr package](#) which allows to fetch the observation directly from FRED database using the function `fredr`.

Here is an excerpt of the data:

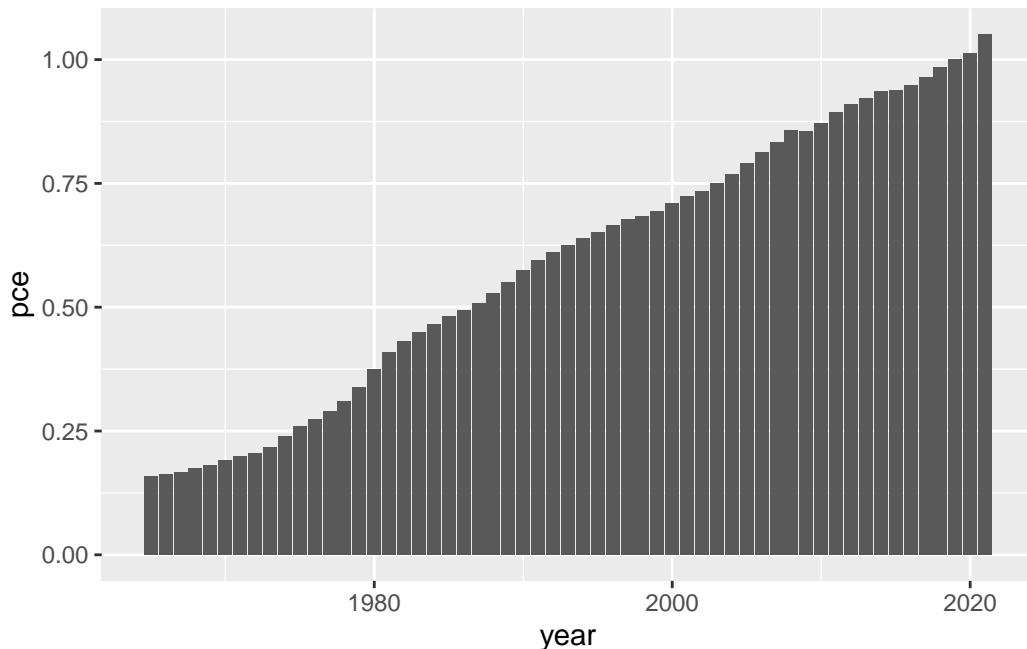
	DATE	PCEPI_NBD20190101
1	1965-01-01	15.92229
2	1966-01-01	16.32477
3	1967-01-01	16.73491
4	1968-01-01	17.38977
5	1969-01-01	18.17313
6	1970-01-01	19.02267

## 9. Collection of exercises

2. Divide the variable PCEPI\_NBD20190101 by 100 and name the resulting variable pce. Additionally, generate a new variable year that contains the respective year. Save the modified dataframe as pce\_c1.
3. Make the following plot:

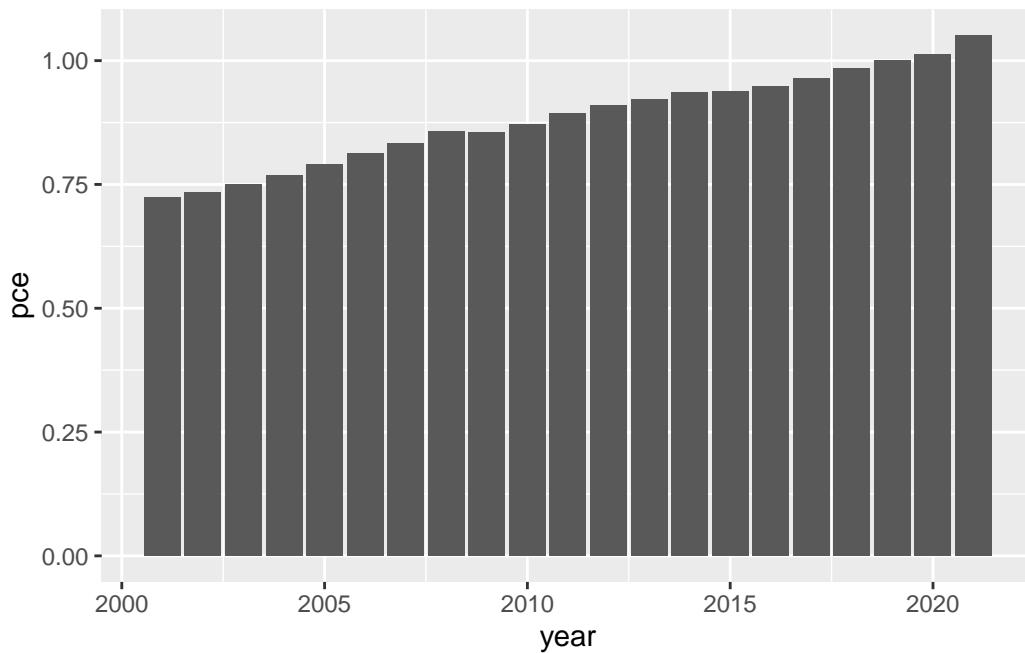


3. Make the following plot:

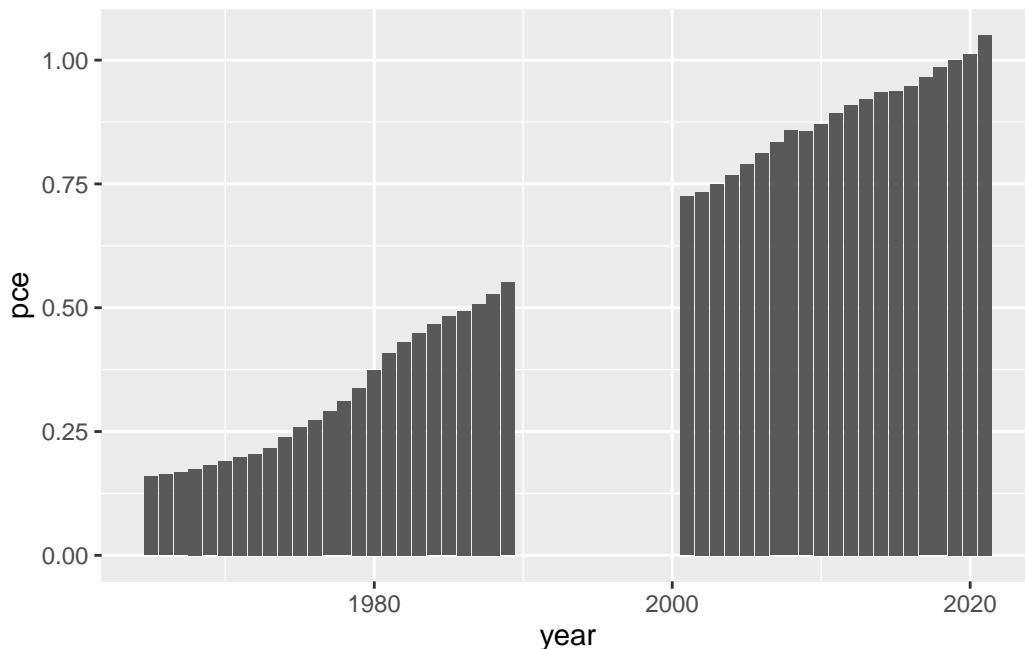


4. Make the following plot:

9. Collection of exercises



5. Make a plot of inflation for all years except the 90s:

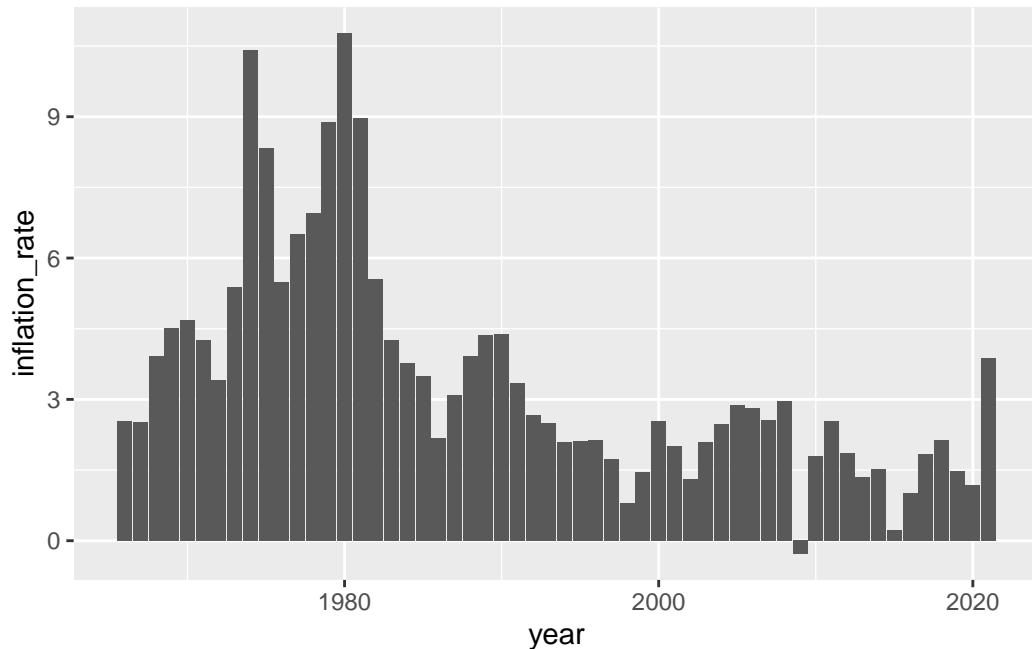


6. Calculate the yearly inflation. Here is an excerpt of how the data should look like:

year	pce	inflation_rate
1	0.1592229	NA
2	0.1632477	2.527777
3	0.1673491	2.512378
4	0.1738977	3.913137
5	0.1817313	4.504717
6	0.1902267	4.674704

7. Plot the yearly inflation rate:

9. Collection of exercises

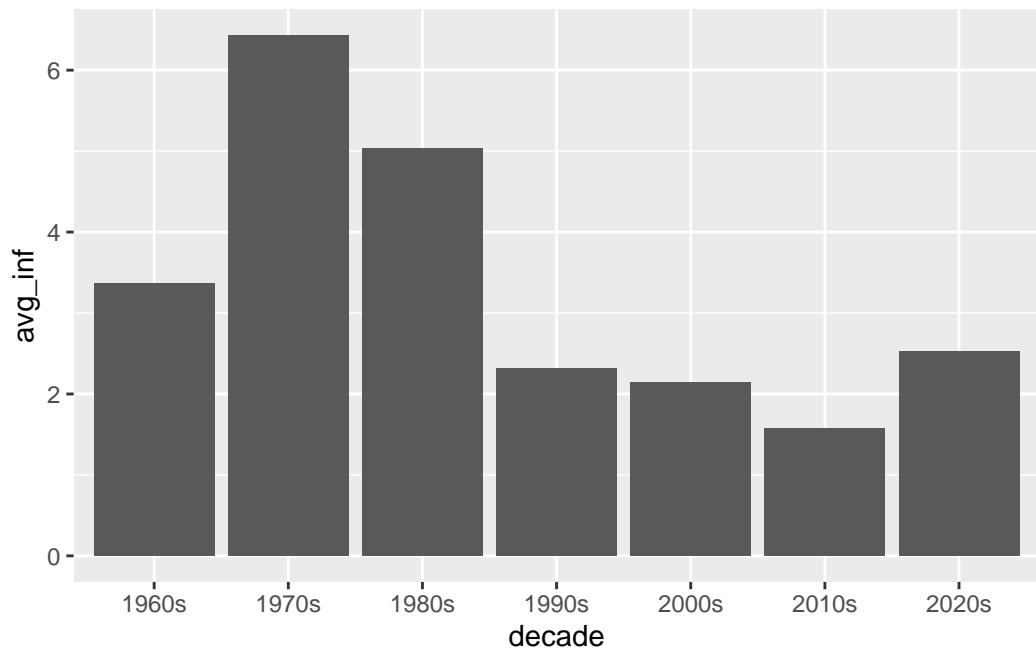


8. Calculate the average inflation rate over all years.
9. Calculate the average inflation rate for each decade:

```
# A tibble: 7 x 2
  decade avg_inf
  <chr>    <dbl>
1 1960s     3.36
2 1970s     6.43
3 1980s     5.03
4 1990s     2.32
5 2000s     2.14
6 2010s     1.57
7 2020s     2.53
```

10. Make the following plot:

## 9. Collection of exercises



### Solution

The script uses the following functions: `aes`, `as.integer`, `case_when`, `filter`, `geom_bar`, `geom_line`, `geom_point`, `ggplot`, `group_by`, `lag`, `mean`, `mutate`, `read.csv`, `select`, `str_sub`, `summarize`.

### R script

## 9. Collection of exercises

```
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, janitor, expss)
# setwd("~/yourdirectory-of-choice")
rm(list = ls())

# read in raw data
# PCEPI <- read.csv("PCEPI.csv")
PCEPI <- read.csv("https://raw.githubusercontent.com/hubchev/courses/main/dta/PCEPI.csv")

# clean data
pce_cl <- PCEPI |>
  mutate(pce = PCEPI_NBD20190101/100) |>
  mutate(year = str_sub(DATE, 1, 4)) |>
  mutate(year = as.integer(year)) |>
  select(pce, year)

# make a plot
pce_cl |>
  ggplot( aes(x=year, y=pce))+
  geom_line() +
  geom_point()

# make a barplot
pce_cl |>
  ggplot( aes(x=year, y=pce))+ 
  geom_bar(stat="identity")

# make a plot for all years from 2000 onwards
pce_cl |>
  filter(year > 2000 ) |>
  ggplot( aes(x=year, y=pce)) +
  geom_bar(stat="identity")

# make a plot for all years except the 90s
pce_cl |>
  filter(year > 2000 | year <1990) |>
  ggplot( aes(x=year, y=pce)) +
  geom_bar(stat="identity")

# calculate yearly inflation
pce_cl <- pce_cl |>
  mutate(inflation_rate = (pce/lag(pce)-1)*100 )

# plot the inflation rate
pce_cl |>
  ggplot( aes(x=year, y=inflation_rate))+ 
  geom_bar(stat="identity")

# what is the average inflation rate
pce_cl |>
  summarize(avg_mpg = mean(inflation_rate, na.rm = TRUE))

# make a variable that indicates the decades 1 for 60s, 2 for 70s, etc.
pce_cl <- pce_cl |>
```

### Output of the R script

```

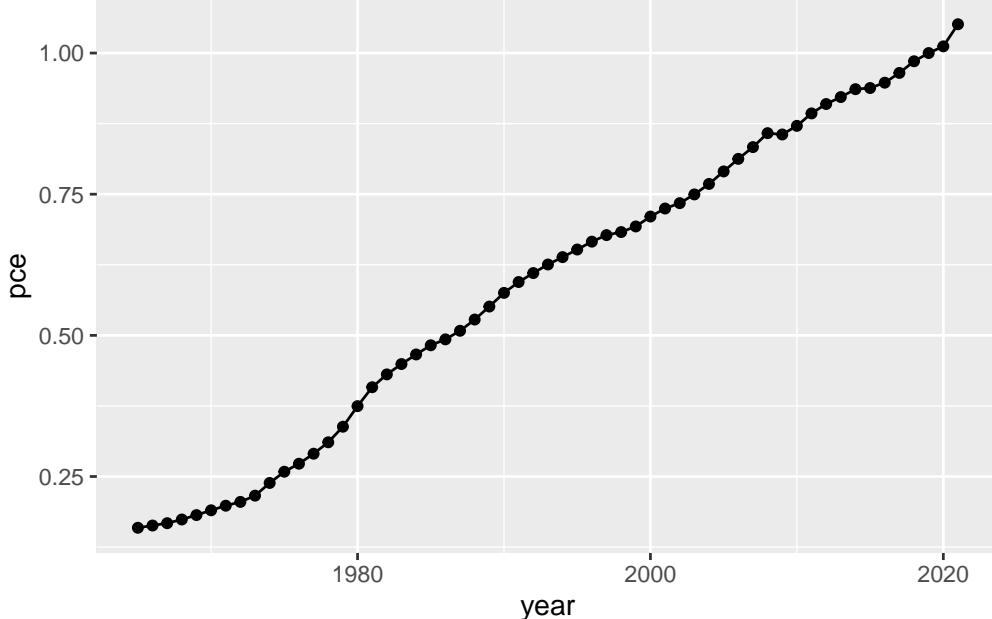
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, janitor, expss)
# setwd("~/yourdirectory-of-choice")
rm(list = ls())


# read in raw data
# PCEPI <- read.csv("PCEPI.csv")
PCEPI <- read.csv("https://raw.githubusercontent.com/hubchev/courses/main/dta/PCEPI.csv")


# clean data
pce_cl <- PCEPI |>
  mutate(pce = PCEPI_NBD20190101/100) |>
  mutate(year = str_sub(DATE, 1, 4)) |>
  mutate(year = as.integer(year)) |>
  select(pce, year)

# make a plot
pce_cl |>
  ggplot( aes(x=year, y=pce))+
  geom_line() +
  geom_point()

```

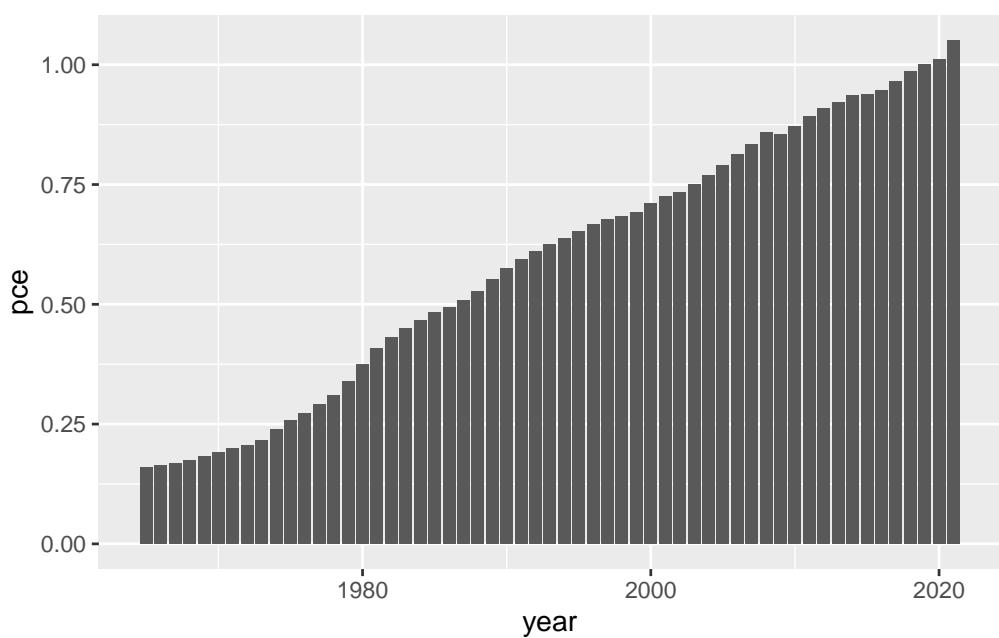


```

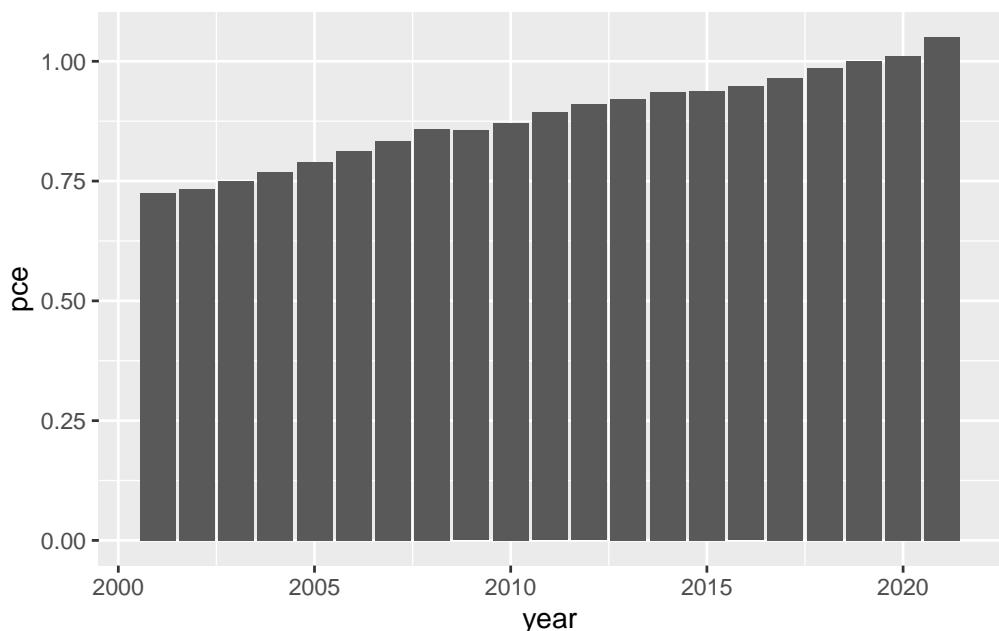
# make a barplot
pce_cl |>
  ggplot( aes(x=year, y=pce))+
  geom_bar(stat="identity")

```

## 9. Collection of exercises

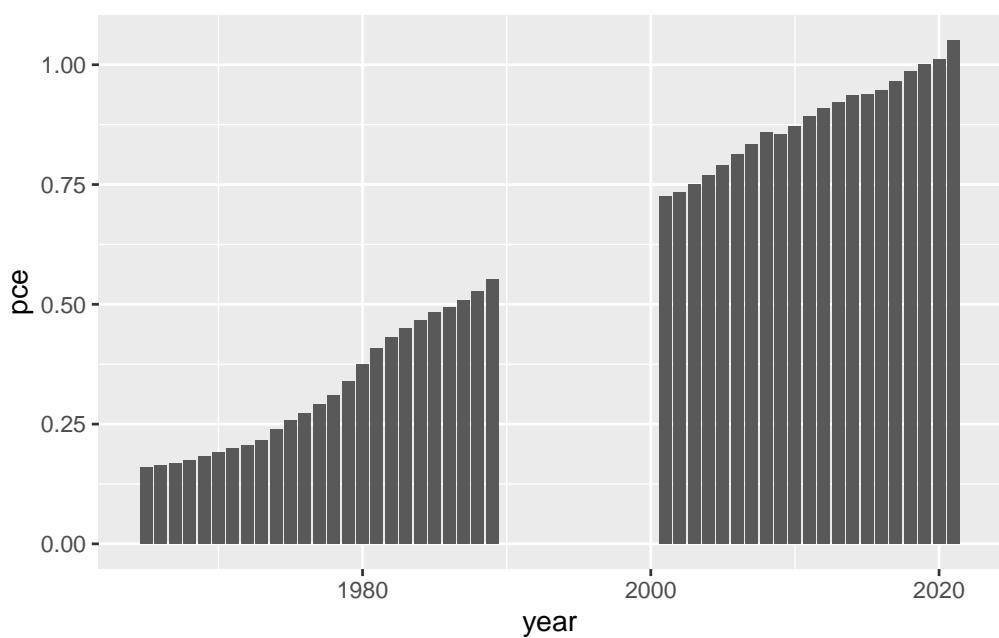


```
# make a plot for all years from 2000 onwards
pce_cl |>
  filter(year > 2000 ) |>
  ggplot( aes(x=year, y=pce)) +
  geom_bar(stat="identity")
```



```
# make a plot for all years except the 90s
pce_cl |>
  filter(year > 2000 | year <1990) |>
  ggplot( aes(x=year, y=pce)) +
  geom_bar(stat="identity")
```

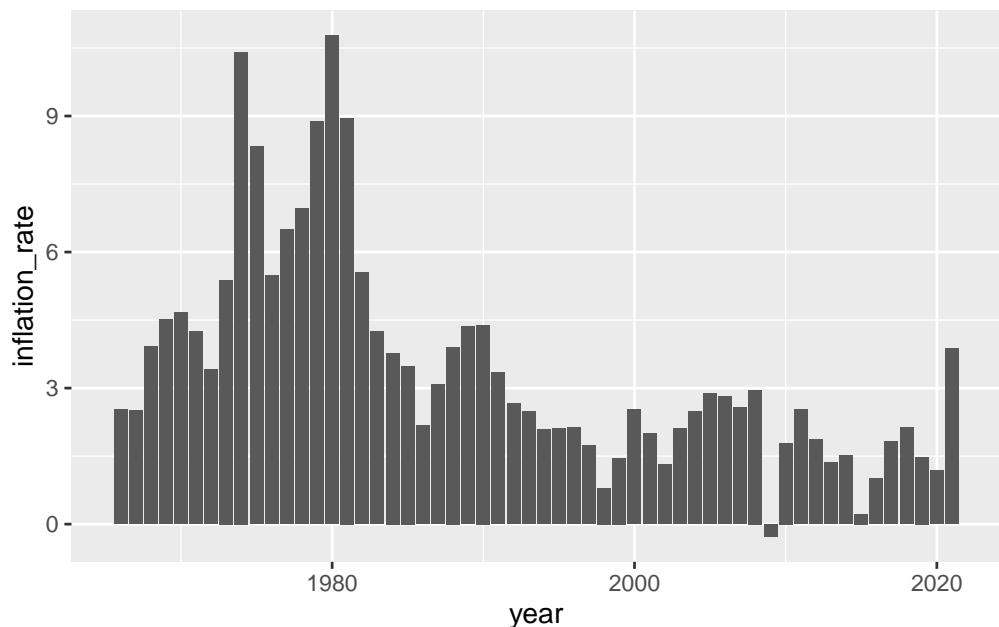
## 9. Collection of exercises



```
# calculate yearly inflation
pce_cl <- pce_cl |>
  mutate(inflation_rate = (pce/lag(pce)-1)*100)

# plot the inflation rate
pce_cl |>
  ggplot( aes(x=year, y=inflation_rate))+
  geom_bar(stat="identity")
```

Warning: Removed 1 row containing missing values or values outside the scale range (`geom\_bar()`).



## 9. Collection of exercises

```
# what is the average inflation rate
pce_cl |>
  summarize(avg_mpg = mean(inflation_rate, na.rm = TRUE))

avg_mpg
1 3.454529

# make a variable that indicates the decades 1 for 60s, 2 for 70s, etc.

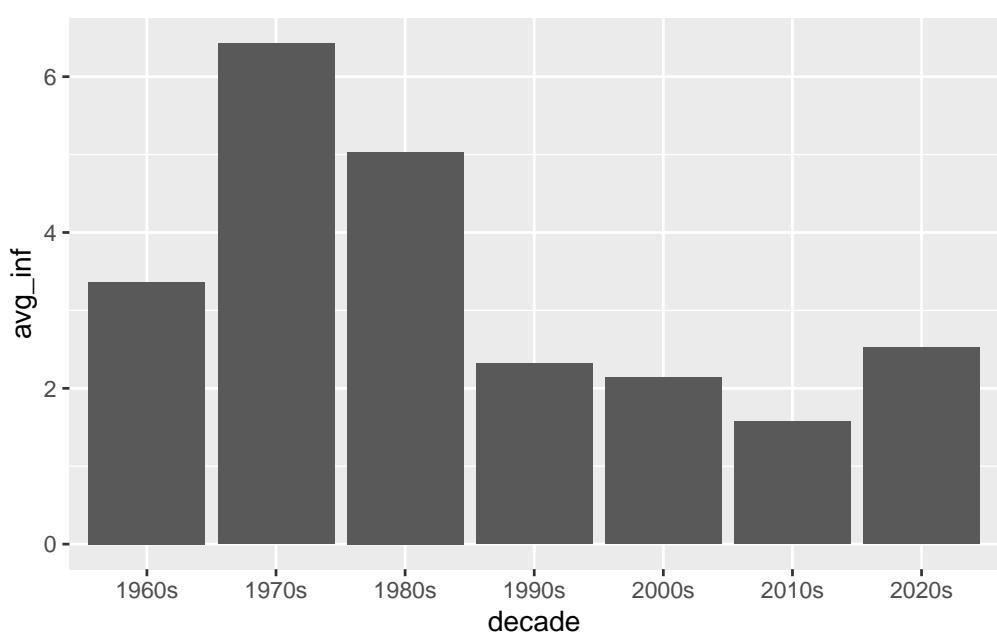
pce_cl <- pce_cl |>
  mutate(decade = case_when(
    year < 1970 ~ "1960s",
    (year > 1969 & year < 1980) ~ "1970s",
    (year > 1979 & year < 1990) ~ "1980s",
    (year > 1989 & year < 2000) ~ "1990s",
    (year > 1999 & year < 2010) ~ "2000s",
    (year > 2009 & year < 2020) ~ "2010s",
    (year > 2019 & year < 2030) ~ "2020s",
  ))

pce_decade <- pce_cl |>
  group_by(decade) |>
  summarize(avg_inf = mean(inflation_rate, na.rm = TRUE))

pce_decade

# A tibble: 7 x 2
  decade avg_inf
  <chr>     <dbl>
1 1960s      3.36
2 1970s      6.43
3 1980s      5.03
4 1990s      2.32
5 2000s      2.14
6 2010s      1.57
7 2020s      2.53

pce_decade |>
  ggplot( aes(x=decade, y=avg_inf))+
  geom_bar(stat="identity")
```



## 9.8. Load the Stata dataset “auto” using R

1. Create a scatter plot illustrating the relationship between the price and weight of a car. Provide a meaningful title for the graph and try to make it clear which car each observation corresponds to.
2. Save this graph in the formats of .png and .pdf.
3. Create a variable `l1p100km` that indicates the fuel consumption of an average car in liters per 100 kilometers. (Note: One gallon is approximately equal to 3.8 liters, and one mile is about 1.6 kilometers.)
4. Create a dummy variable `larger6000` that is equal to 1 if the price of a car is above \$6000.
5. Now, search for the “most unreasonable poser car” that costs no more than \$6000. A *poser* car is defined as one that is expensive, has a large turning radius, consumes a lot of fuel, and is often defective (`rep78` is low). For this purpose, create a metric indicator for each corresponding variable that indicates a value of 1 for the car that is the most unreasonable in that variable and 0 for the most reasonable car. All other cars should fall between 0 and 1.

### Solution

The script uses the following functions: `aes`, `arrange`, `desc`, `dir.create`, `dir.exists`, `filter`, `geom_point`, `geom_text_repel`, `ggplot`, `head`, `ifelse`, `max`, `min`, `min_max_norm`, `mutate`, `na.omit`, `read_dta`, `select`, `tail`, `theme_minimal`, `xlab`, `ylab`.

### R script

## 9. Collection of exercises

```
# Load the required libraries
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, haven, ggrepel)

# setwd("~/Dropbox/hsf/23-ws/R_begin")

rm(list = ls())

# Read the Stata dataset
auto <- read_dta("http://www.stata-press.com/data/r8/auto.dta")

# Create a scatter plot of price vs. weight
scatter_plot <- ggplot(auto, aes(x = mpg, y = price, label = make)) +
  geom_point() +
  geom_text_repel() +
  xlab("Miles per Gallon") +
  ylab("Price in Dollar") +
  theme_minimal()

scatter_plot

# Create "fig" directory if it doesn't already exist
if (!dir.exists("fig")) {
  dir.create("fig")
}

# Save the scatter plot in different formats
# ggsave("fig/scatter_plot.png", plot = scatter_plot, device = "png")
# ggsave("fig/scatter_plot.pdf", plot = scatter_plot, device = "pdf")

# Create 'lp100km' variable for fuel consumption
n_auto <- auto |>
  mutate(lp100km = (1 / (mpg * 1.6 / 3.8)) * 100)

# Create 'larger6000' dummy variable
n_auto <- n_auto |>
  mutate(larger6000 = ifelse(price > 6000, 1, 0))

# Normalize variables

## Do it slowly
n_auto <- n_auto |>
  mutate(sprice = (price - min(auto$price)) / (max(auto$price) - min(auto$price)))

n_auto <- n_auto |>
  filter(larger6000 == 0)

## Do it with a self-written function
min_max_norm <- function(x) {
  (x - min(x, na.rm = TRUE)) / (max(x, na.rm = TRUE) - min(x, na.rm = TRUE))
}
```

### Output of the R script

```
# Load the required libraries
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, haven, ggrepel)

# setwd("~/Dropbox/hsf/23-ws/R_begin")

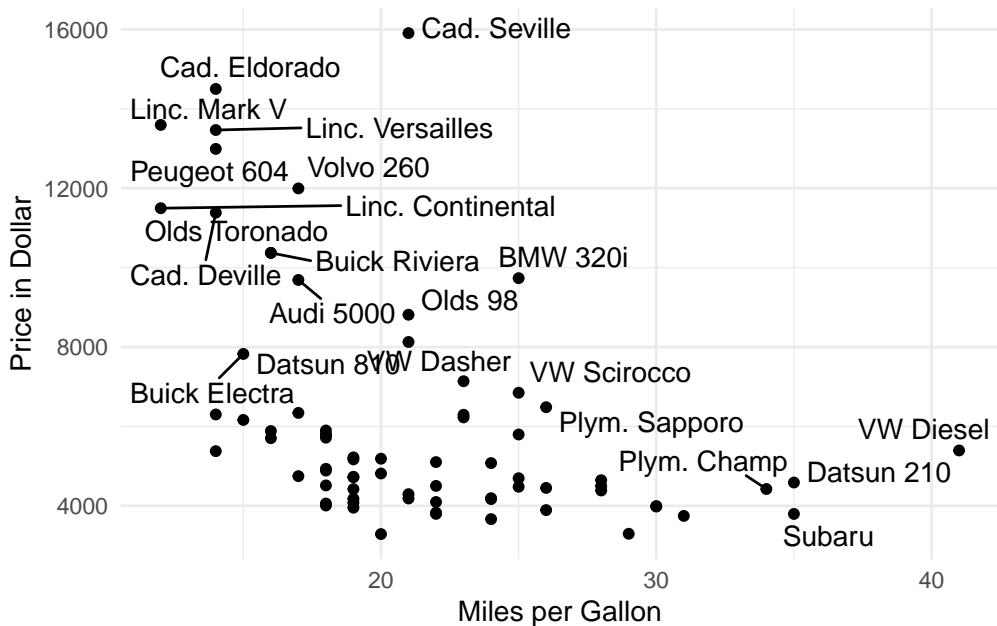
rm(list = ls())

# Read the Stata dataset
auto <- read_dta("http://www.stata-press.com/data/r8/auto.dta")

# Create a scatter plot of price vs. weight
scatter_plot <- ggplot(auto, aes(x = mpg, y = price, label = make)) +
  geom_point() +
  geom_text_repel() +
  xlab("Miles per Gallon") +
  ylab("Price in Dollar") +
  theme_minimal()

scatter_plot
```

Warning: ggrepel: 52 unlabeled data points (too many overlaps). Consider increasing max.overlaps



## 9. Collection of exercises

```
# Create "fig" directory if it doesn't already exist
if (!dir.exists("fig")) {
  dir.create("fig")
}

# Save the scatter plot in different formats
# ggsave("fig/scatter_plot.png", plot = scatter_plot, device = "png")
# ggsave("fig/scatter_plot.pdf", plot = scatter_plot, device = "pdf")

# Create 'lp100km' variable for fuel consumption
n_auto <- auto |>
  mutate(lp100km = (1 / (mpg * 1.6 / 3.8)) * 100)

# Create 'larger6000' dummy variable
n_auto <- n_auto |>
  mutate(larger6000 = ifelse(price > 6000, 1, 0))

# Normalize variables

## Do it slowly
n_auto <- n_auto |>
  mutate(sprice = (price - min(auto$price)) / (max(auto$price) - min(auto$price)))

n_auto <- n_auto |>
  filter(larger6000 == 0)

## Do it with a self-written function
min_max_norm <- function(x) {
  (x - min(x, na.rm = TRUE)) / (max(x, na.rm = TRUE) - min(x, na.rm = TRUE))
}

n_auto <- n_auto |>
  mutate(smpg = min_max_norm(mpg)) |>
  mutate(sturn = min_max_norm(turn)) |>
  mutate(slp100km = min_max_norm(lp100km)) |>
  mutate(sprice = min_max_norm(price)) |>
  mutate(srep78 = min_max_norm(rep78))

## With a loop:

# vars_to_normalize <- c("mpg", "turn", "lp100km", "price", "rep78")
#
# # Loop through the selected variables and apply min_max_norm
# for (var in c("mpg", "turn", "lp100km", "price", "rep78")) {
#   auto <- auto |>
#     mutate (!!paste0("s", var) := min_max_norm (!!sym(var))) |>
#     select(make, starts_with("s"))
# }

## mpg and rep78 need to be changed because a SMALL value is poser-like
n_auto <- n_auto |>
  mutate(smpg = 1 - smpg) |>
  mutate(srep78 = 1 - srep78)  141
```

## 9. Collection of exercises

```
# A tibble: 5 x 2
  make          poser
  <chr>        <dbl>
1 Dodge Magnum 0.888
2 Pont. Firebird 0.782
3 Merc. Cougar 0.763
4 Buick LeSabre 0.754
5 Pont. Grand Prix 0.720
```

```
df_poser <- n_auto |>
  filter(larger6000 == 0) |>
  arrange(desc(poser)) |>
  select(make, poser) |>
  na.omit()

# Five top poser cars
head(df_poser, 15)
```

```
# A tibble: 15 x 2
  make          poser
  <chr>        <dbl>
1 Dodge Magnum 0.888
2 Pont. Firebird 0.782
3 Merc. Cougar 0.763
4 Buick LeSabre 0.754
5 Pont. Grand Prix 0.720
6 Chev. Impala 0.702
7 Dodge Diplomat 0.690
8 Chev. Monte Carlo 0.684
9 Pont. Catalina 0.678
10 Olds Cutl Supr 0.671
11 Plym. Volare 0.665
12 Buick Regal 0.663
13 Olds Cutlass 0.629
14 Olds Starfire 0.626
15 AMC Pacer 0.619
```

```
# Five top non-poser cars
tail(df_poser, 5)
```

```
# A tibble: 5 x 2
  make          poser
  <chr>        <dbl>
1 VW Diesel    0.261
2 Dodge Colt   0.227
3 Toyota Corolla 0.195
4 Datsun 210   0.195
5 Subaru       0.178
```

```
# unload packages
suppressMessages(pacman::p_unload(tidyverse, haven, ggrepel))
```

## 9.9. DatasauRus

Figure 9.1.: The logo of the DatasauRus package



Source: <https://github.com/jumpingrivers/datasauRus> —

- a) Load the packages `datasauRus` and `tidyverse`. If necessary, install these packages.
- b) The packaged `datasauRus` comes with a dataset in two different formats: `datasaurus_dozen` and `datasaurus_dozen_wide`. Store them as `ds` and `ds_wide`.
- c) Open and read the R vignette of the `datasauRus` package. Also open the R documentation of the dataset `datasaurus_dozen`.
- d) Explore the dataset: What are the dimensions of this dataset? Look at the descriptive statistics.
- e) How many unique values does the variable `dataset` of the tibble `ds` have? Hint: The function `unique()` return the unique values of a variable and the function `length()` returns the length of a vector, such as the unique elements.
- f) Compute the mean values of the `x` and `y` variables for each entry in `dataset`. Hint: Use the `group_by()` function to group the data by the appropriate column and then the `summarise()` function to calculate the mean.
- g) Compute the standard deviation, the correlation, and the median in the same way. Round the numbers.
- h) What can you conclude?
- i) Plot all datasets of `ds`. Hide the legend. Hint: Use the `facet_wrap()` and the `theme()` function.
- j) Create a loop that generates separate scatter plots for each unique datatset of the tibble `ds`. Export each graph as a png file.
- k) Watch the video [Animating the Datasaurus Dozen Dataset in R](#) from The Data Digest on YouTube.

 Solution

The script uses the following functions: `aes`, `cor`, `dim`, `dir.create`, `dir.exists`, `facet_wrap`, `filter`, `geom_point`, `ggplot`, `ggsave`, `glimpse`, `group_by`, `head`, `labs`, `length`, `mean`, `median`, `paste`, `paste0`, `round`, `sd`, `select`, `summarise`, `summary`, `theme`, `theme_bw`, `unique`, `view`.

**R script**

## 9. Collection of exercises

```
# setwd("/home/sthu/Dropbox/hsf/23-ws/ds_mim/")
rm(list = ls())
```

```
# Load the packages datasauRus and tidyverse. If necessary, install these packages.
```

```
# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(datasauRus, tidyverse)
```

```
# The packagedatasauRus comes with a dataset in two different formats:
#   datasaurus_dozen and datasaurus_dozen_wide. Store them as ds and ds_wide.
```

```
ds <- datasaurus_dozen
ds_wide <- datasaurus_dozen_wide
```

```
# Open and read the R vignette of the datasauRus package.
```

```
#   Also open the R documentation of the dataset datasaurus_dozen.
```

```
??datasaurus
```

```
# Explore the dataset: What are the dimensions of this dataset? Look at the descriptive
```

```
ds
dim(ds)
head(ds)
glimpse(ds)
view(ds)
summary(ds)
```

```
# How many unique values does the variable dataset of the tibble ds have?
```

```
#   Hint: The function unique() return the unique values of a variable and the
#         function length() returns the length of a vector, such as the unique elements.
```

```
unique(ds$dataset)
```

```
unique(ds$dataset) |>
  length()
```

```
# Compute the mean values of the x and y variables for each entry in dataset.
```

```
#   Hint: Use the group_by() function to group the data by the appropriate column and
#         then the summarise() function to calculate the mean.
```

```
ds |>
  group_by(dataset) |>
  summarise(
    mean_x = mean(x),
    mean_y = mean(y)
  )
```

```
# Compute the standard deviation, the correlation, and the median in the same way. Round
```

```
ds |>
  group_by(dataset) |>
  summarise(
    mean_x = round(mean(x), 2), 145
    mean_y = round(mean(y), 2),
```

## Output of the R script

```
# setwd("/home/sthu/Dropbox/hsf/23-ws/ds_mim/")
rm(list = ls())

# Load the packages datasauRus and tidyverse. If necessary, install these packages.

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(datasauRus, tidyverse)

# The packagedatasauRus comes with a dataset in two different formats:
#   datasaurus_dozen and datasaurus_dozen_wide. Store them as ds and ds_wide.

ds <- datasaurus_dozen
ds_wide <- datasaurus_dozen_wide

# Open and read the R vignette of the datasauRus package.
#   Also open the R documentation of the dataset datasaurus_dozen.

??datasaurus

# Explore the dataset: What are the dimensions of this dataset? Look at the descriptive statistics of ds.

# A tibble: 1,846 x 3
  dataset      x      y
  <chr>    <dbl> <dbl>
1 dino      55.4  97.2
2 dino      51.5  96.0
3 dino      46.2  94.5
4 dino      42.8  91.4
5 dino      40.8  88.3
6 dino      38.7  84.9
7 dino      35.6  79.9
8 dino      33.1  77.6
9 dino      29.0  74.5
10 dino     26.2  71.4
# i 1,836 more rows

dim(ds)

[1] 1846     3

head(ds)

# A tibble: 6 x 3
  dataset      x      y
  <chr>    <dbl> <dbl>
1 dino      55.4  97.2
2 dino      51.5  96.0
3 dino      46.2  94.5
4 dino      42.8  91.4
5 dino      40.8  88.3
6 dino      38.7  84.9
```

## 9. Collection of exercises

```
<chr> <dbl> <dbl>
1 dino    55.4  97.2
2 dino    51.5  96.0
3 dino    46.2  94.5
4 dino    42.8  91.4
5 dino    40.8  88.3
6 dino    38.7  84.9
```

```
glimpse(ds)
```

```
Rows: 1,846
Columns: 3
$ dataset <chr> "dino", "dino", "dino", "dino", "dino", "dino", "dino"~
$ x         <dbl> 55.3846, 51.5385, 46.1538, 42.8205, 40.7692, 38.7179, 35.6410, ~
$ y         <dbl> 97.1795, 96.0256, 94.4872, 91.4103, 88.3333, 84.8718, 79.8718, ~
```

```
view(ds)
summary(ds)
```

	x	y
Length:1846	Min. :15.56	Min. : 0.01512
Class :character	1st Qu.:41.07	1st Qu.:22.56107
Mode :character	Median :52.59	Median :47.59445
	Mean :54.27	Mean :47.83510
	3rd Qu.:67.28	3rd Qu.:71.81078
	Max. :98.29	Max. :99.69468

```
# How many unique values does the variable dataset of the tibble ds have?
# Hint: The function unique() return the unique values of a variable and the
# function length() returns the length of a vector, such as the unique elements.
```

```
unique(ds$dataset)
```

```
[1] "dino"      "away"       "h_lines"    "v_lines"    "x_shape"
[6] "star"       "high_lines" "dots"       "circle"     "bullseye"
[11] "slant_up"   "slant_down" "wide_lines"
```

```
unique(ds$dataset) |>
  length()
```

```
[1] 13
```

## 9. Collection of exercises

```
# Compute the mean values of the x and y variables for each entry in dataset.  
# Hint: Use the group_by() function to group the data by the appropriate column and  
# then the summarise() function to calculate the mean.  
  
ds |>  
  group_by(dataset) |>  
  summarise(  
    mean_x = mean(x),  
    mean_y = mean(y)  
  )  
  
# A tibble: 13 x 3  
  dataset   mean_x  mean_y  
  <chr>     <dbl>   <dbl>  
1 away      54.3    47.8  
2 bullseye  54.3    47.8  
3 circle    54.3    47.8  
4 dino      54.3    47.8  
5 dots      54.3    47.8  
6 h_lines   54.3    47.8  
7 high_lines 54.3    47.8  
8 slant_down 54.3    47.8  
9 slant_up   54.3    47.8  
10 star     54.3    47.8  
11 v_lines   54.3    47.8  
12 wide_lines 54.3    47.8  
13 x_shape  54.3    47.8  
  
# Compute the standard deviation, the correlation, and the median in the same way. R  
  
ds |>  
  group_by(dataset) |>  
  summarise(  
    mean_x = round(mean(x), 2),  
    mean_y = round(mean(y), 2),  
    sd_x = round(sd(x), 2),  
    sd_y = round(sd(y), 2),  
    med_x = round(median(x), 2),  
    med_y = round(median(y), 2),  
    cor = round(cor(x, y), digits = 4)  
  )  
  
# A tibble: 13 x 8  
  dataset   mean_x  mean_y  sd_x  sd_y med_x  med_y     cor  
  <chr>     <dbl>   <dbl>  <dbl> <dbl> <dbl>   <dbl>    <dbl>  
1 away      54.3    47.8   16.8  26.9  53.3   47.5 -0.0641  
2 bullseye  54.3    47.8   16.8  26.9  53.8   47.4 -0.0686  
3 circle    54.3    47.8   16.8  26.9  54.0   51.0 -0.0683  
4 dino      54.3    47.8   16.8  26.9  53.3   46.0 -0.0645
```

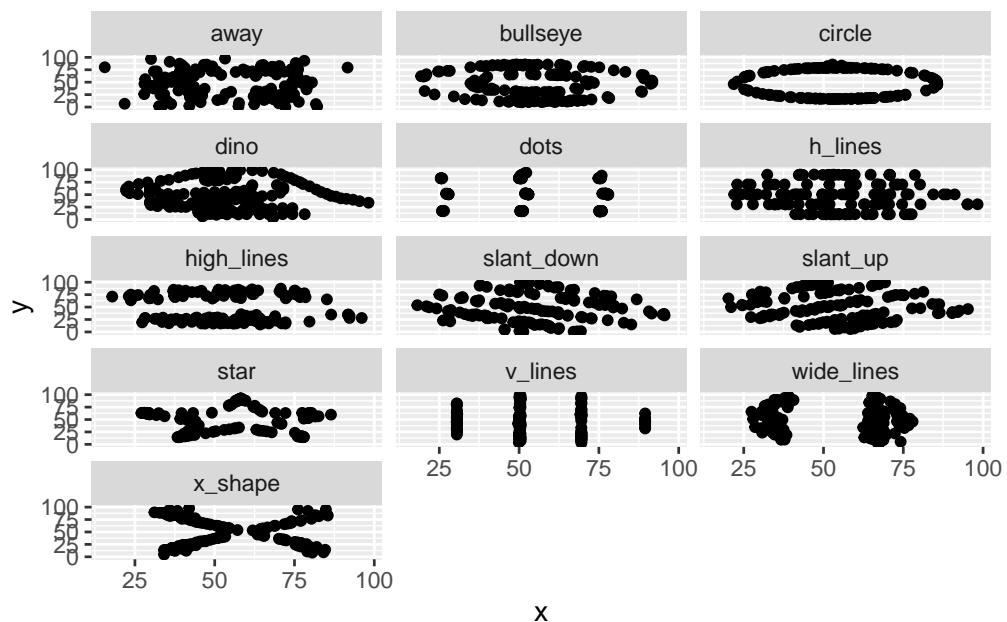
## 9. Collection of exercises

5 dots	54.3	47.8	16.8	26.9	51.0	51.3	-0.0603
6 h_lines	54.3	47.8	16.8	26.9	53.1	50.5	-0.0617
7 high_lines	54.3	47.8	16.8	26.9	54.2	32.5	-0.0685
8 slant_down	54.3	47.8	16.8	26.9	53.1	46.4	-0.069
9 slant_up	54.3	47.8	16.8	26.9	54.3	45.3	-0.0686
10 star	54.3	47.8	16.8	26.9	56.5	50.1	-0.063
11 v_lines	54.3	47.8	16.8	26.9	50.4	47.1	-0.0694
12 wide_lines	54.3	47.8	16.8	26.9	64.6	46.3	-0.0666
13 x_shape	54.3	47.8	16.8	26.9	47.1	39.9	-0.0656

```
# What can you conclude?
# --> The standard deviation, the mean, and the correlation are basically the
# same for all datasets. The median is different.

# Plot all datasets of ds. Hide the legend. Hint: Use the facet_wrap() and the theme()

ggplot(ds, aes(x = x, y = y)) +
  geom_point() +
  facet_wrap(~dataset, ncol = 3) +
  theme(legend.position = "none")
```



## 9. Collection of exercises

```
# Create a loop that generates separate scatter plots for each unique dataset of the
#   Export each graph as a png file.

# Assuming uni_ds is a vector of unique values for the 'dataset' variable
uni_ds <- unique(ds$dataset)

# Create the 'pic' folder if it doesn't exist
if (!dir.exists("fig")) {
  dir.create("fig")
}

for (uni_v in uni_ds) {
  # Select data for the current value
  subset_ds <- ds |>
    filter(dataset == uni_v) |>
    select(x, y)

  # Make plot
  graph <- ggplot(subset_ds, aes(x = x, y = y)) +
    geom_point() +
    labs(
      title = paste("Dataset:", uni_v),
      x = "X",
      y = "Y"
    ) +
    theme_bw()

  # Save the plot as a PNG file
  filename <- paste0("fig/", "plot_ds_", uni_v, ".png")
  ggsave(filename, plot = graph)
}
```

```
Saving 5.5 x 3.5 in image
```

```
# unload packages
suppressMessages(pacman::p_unload(datasauRus, tidyverse))
```

## 9.10. Convergence

The dataset convergence.dta, see <https://github.com/hubchev/courses/blob/main/dta/convergence.dta>, contains the per capita GDP of 1960 (`gdppc60`) and the average growth rate of GDP per capita between 1960 and 1995 (`growth`) for different countries (`country`), as well as 3 dummy variables indicating the belonging of a country to the region Asia (`asia`), Western Europe (`weurope`) or Africa (`africa`).

- Some countries are not assigned to a certain country group. Name the countries which are assigned to be part of Western Europe, Africa or Asia. If you find countries that are members of the EU, assign them a ‘1’ in the variable `weurope`.
- Create a table that shows the average GDP per capita for all available points in time. Group by Western European, Asian, African, and the remaining countries.
- Create the growth rate of GDP per capita from 1960 to 1995 and call it `gdpgrowth`. (Note: The log value X minus the log value X of the previous period is approximately equal to the growth rate).
- Calculate the unconditional convergence of all countries by constructing a graph in which a scatterplot shows the GDP per capita growth rate between 1960 and 1995 (`gdpgrowth`) on the y-axis and the 1960 GDP per capita (`gdppc60`) on the x-axis. Add to the same graph the estimated linear relationship. You do not need to label the graph further, just two things: title the graph `world` and label the individual observations with the country names.
- Create three graphs describing the same relationship for the sample of Western European, African and Asian countries. Title the graph accordingly with `weurope`, `africa` and `asia`.
- Combine the four graphs into one image. Discuss how an upward or downward sloping regression line can be interpreted.
- Estimate the relationships illustrated in the 4 graphs using the least squares method. Present the 4 estimation results in a table, indicating the significance level with stars. In addition, the Akaike information criterion, and the number of observations should be displayed in the table. Interpret the four estimation results regarding their significance.
- Put the data set into the so-called long format and calculate the GDP per capita growth rates for the available time points in the countries.

### Solution

The script uses the following functions: `aes`, `as.numeric`, `c`, `cor`, `describe`, `diff`, `filter`, `gather`, `geom_point`, `geom_text`, `ggarrange`, `ggplot`, `ggtitle`, `group_by`, `head`, `ifelse`, `lag`, `list`, `lm`, `log`, `mean`, `mutate`, `names`, `read_dta`, `select`, `set_label`, `starts_with`, `stat_smooth`, `stat_desc`, `str`, `subset`, `substr`, `summarise`, `summarise_all`, `summarise_at`, `summary`, `tab_model`, `tail`, `vars`, `view`.

### R script

## 9. Collection of exercises

```
# Convergence

# set working directory
# setwd("/home/sthu/Dropbox/hsf/github/courses/")

# clear the environment
rm(list = ls())

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(
  haven, tidyverse, vtable, gtsummary, pastecs, Hmisc,
  sjlabelled, tis, ggpibr, sjPlot, psych
)

# import data
data <- read_dta("https://github.com/hubchev/courses/raw/main/dta/convergence.dta")

# inspect data
names(data)
str(data)
data
head(data)
tail(data)
summary(data)
view(data)

# library(vtable)
# vtable(data, missing=TRUE)

# library(pastecs)
stat.desc(data)

# library(Hmisc)
describe(data)

# library(gtsummary)
#tbl_summary(data)

# check the assignments of countries to continents
data |>
  select(country, africa, asia, weurope) |>
  view()

data <- mutate(data, x_1 = africa + asia + weurope)

data |>
  filter(x_1 == 0) |>
  select(africa, asia, weurope, country) |>
  view()

# correct the assignment manually
data$weurope[data$country == "Austria"] <- 1
data$weurope[data$country == "Greece"] <- 1
data$weurope[data$country == "Cyprus"] <- 152
```

## Output of the R script

```
# Convergence

# set working directory
# setwd("/home/sthu/Dropbox/hsf/github/courses/")

# clear the environment
rm(list = ls())

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(
  haven, tidyverse, vtable, gtsummary, pastecs, Hmisc,
  sjlabelled, tis, ggpubr, sjPlot, psych
)

# import data
data <- read_dta("https://github.com/hubchev/courses/raw/main/dta/convergence.dta")

# inspect data
names(data)
```

[1] "country" "gdppc60" "gdppc65" "gdppc70" "gdppc75" "gdppc80" "gdppc85"  
[8] "gdppc90" "gdppc95" "africa" "asia" "weurope" "growth"

```
str(data)
```

```
tibble [107 x 13] (S3:tbl_df/tbl/data.frame)
$ country: chr [1:107] "Algeria" "Angola" "Argentina" "Australia" ...
..- attr(*, "format.stata")= chr "%24s"
$ gdppc60: num [1:107] 2848 2642 7879 11436 7842 ...
..- attr(*, "label")= chr "real gdp per capita 1960"
..- attr(*, "format.stata")= chr "%9.0g"
$ gdppc65: num [1:107] 3536 3072 8802 13192 9387 ...
..- attr(*, "label")= chr "real gdp per capita 1965"
..- attr(*, "format.stata")= chr "%9.0g"
$ gdppc70: num [1:107] 3670 3558 9903 15842 11946 ...
..- attr(*, "label")= chr "real gdp per capita 1970"
..- attr(*, "format.stata")= chr "%9.0g"
$ gdppc75: num [1:107] 3917 2230 10609 16716 14198 ...
..- attr(*, "label")= chr "real gdp per capita 1975"
..- attr(*, "format.stata")= chr "%9.0g"
$ gdppc80: num [1:107] 5094 2059 11359 18300 16869 ...
..- attr(*, "label")= chr "real gdp per capita 1980"
..- attr(*, "format.stata")= chr "%9.0g"
$ gdppc85: num [1:107] 5876 1988 9246 19669 17919 ...
..- attr(*, "label")= chr "real gdp per capita 1985"
..- attr(*, "format.stata")= chr "%9.0g"
```

## 9. Collection of exercises

```
$ gdppc90: num [1:107] 5307 2081 7716 21446 21178 ...
..- attr(*, "label")= chr "real gdp per capita 1990"
..- attr(*, "format.stata")= chr "%9.0g"
$ gdppc95: num [1:107] 4935 1339 10973 23827 22474 ...
..- attr(*, "label")= chr "real gdp per capita 1995"
..- attr(*, "format.stata")= chr "%9.0g"
$ africa : num [1:107] 1 1 0 0 0 0 0 1 0 ...
..- attr(*, "label")= chr "=1 if in Africa"
..- attr(*, "format.stata")= chr "%8.0g"
$ asia : num [1:107] 0 0 0 0 0 1 0 0 0 ...
..- attr(*, "label")= chr "=1 if in Asia"
..- attr(*, "format.stata")= chr "%8.0g"
$ weurope: num [1:107] 0 0 0 0 0 0 0 1 0 0 ...
..- attr(*, "label")= chr "=1 if in Western Europe"
..- attr(*, "format.stata")= chr "%8.0g"
$ growth : num [1:107] 0.55 -0.68 0.331 0.734 1.053 ...
..- attr(*, "format.stata")= chr "%9.0g"
```

`data`

```
# A tibble: 107 x 13
  country   gdppc60 gdppc65 gdppc70 gdppc75 gdppc80 gdppc85 gdppc90 gdppc95
  <chr>     <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
1 Algeria    2848.    3536.    3670.    3917.    5094.    5876.    5307.    4935.
2 Angola     2642.    3072.    3558.    2230.    2059.    1988.    2081.    1339.
3 Argentina   7879.    8802.    9903.   10609.   11359.   9246.    7716.    10973.
4 Australia   11436.   13192.   15842.   16716.   18300.   19669.   21446.   23827.
5 Austria     7842.    9387.   11946.   14198.   16869.   17919.   21178.   22474.
6 Bangladesh  1130.    1164.    1181.    1030.    1040.    1245.    1366.    1568.
7 Barbados    3632.    4632.    6456.    8827.    10911.   11090.   14411.   14636.
8 Belgium     8314.    10454.   12980.   15024.   17451.   18109.   21246.   22356.
9 Benin       1140.    1188.    1170.    1048.    1069.    1252.    1069.    1139.
10 Bolivia    2516.    2880.    2670.    3124.    3264.    2718.    2615.    2795.
# i 97 more rows
# i 4 more variables: africa <dbl>, asia <dbl>, weurope <dbl>, growth <dbl>
```

`head(data)`

```
# A tibble: 6 x 13
  country   gdppc60 gdppc65 gdppc70 gdppc75 gdppc80 gdppc85 gdppc90 gdppc95 africa
  <chr>     <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
1 Algeria    2848.    3536.    3670.    3917.    5094.    5876.    5307.    4935.    1
2 Angola     2642.    3072.    3558.    2230.    2059.    1988.    2081.    1339.    1
3 Argent~    7879.    8802.    9903.   10609.   11359.   9246.    7716.    10973.   0
4 Austra~   11436.   13192.   15842.   16716.   18300.   19669.   21446.   23827.   0
5 Austria     7842.    9387.   11946.   14198.   16869.   17919.   21178.   22474.   0
6 Bangla~    1130.    1164.    1181.    1030.    1040.    1245.    1366.    1568.   0
# i 3 more variables: asia <dbl>, weurope <dbl>, growth <dbl>
```

`tail(data)`

## 9. Collection of exercises

```
# A tibble: 6 x 13
  country gdppc60 gdppc65 gdppc70 gdppc75 gdppc80 gdppc85 gdppc90 gdppc95 africa
  <chr>    <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
1 United~  10341.  11633.  12917.  14072.  15302.  16878.  19585.  20963.  0
2 United~  13118.  15697.  17478.  19284.  22806.  25251.  28281.  30366.  0
3 Uruguay  6279.   5936.   6553.   6949.   8580.   6625.   7763.   9399.  0
4 Venezu~  8381.   10618.  11253.  8815.   8516.   7274.   7431.   7582.  0
5 Zambia   1290.   1564.   1427.   1446.   1324.   1167.   1091.   870.   1
6 Zimbab~  1317.   1539.   2303.   2694.   2816.   2923.   3115.   2832.  1
# i 3 more variables: asia <dbl>, weurope <dbl>, growth <dbl>
```

```
summary(data)
```

	country	gdppc60	gdppc65	gdppc70	gdppc75	gdppc80	gdppc85	gdppc90	africa
Length:	107	Min. : 407.8	Min. : 513.6	Min. : 354.5					
Class :	character	1st Qu.: 1153.2	1st Qu.: 1364.5	1st Qu.: 1487.9					
Mode :	character	Median : 2484.7	Median : 2884.4	Median : 3072.2					
		Mean : 3634.3	Mean : 4367.5	Mean : 5128.4					
		3rd Qu.: 4354.0	3rd Qu.: 5873.3	3rd Qu.: 6994.6					
		Max. :16010.3	Max. :18928.9	Max. :22030.9					
	gdppc75	gdppc80	gdppc85	gdppc90					
	Min. : 617.9	Min. : 473.6	Min. : 542.3	Min. : 527.7					
	1st Qu.: 1480.7	1st Qu.: 1708.6	1st Qu.: 1598.8	1st Qu.: 1829.0					
	Median : 3741.7	Median : 4306.2	Median : 4200.7	Median : 4034.0					
	Mean : 5759.1	Mean : 6553.6	Mean : 6900.3	Mean : 7775.1					
	3rd Qu.: 8355.8	3rd Qu.: 9968.6	3rd Qu.:10037.2	3rd Qu.:11716.2					
	Max. :21808.9	Max. :23860.1	Max. :25251.4	Max. :28744.1					
	gdppc95	africa	asia	weurope					
	Min. : 499.3	Min. :0.0000	Min. :0.0000	Min. :0.0000					
	1st Qu.: 1673.7	1st Qu.:0.0000	1st Qu.:0.0000	1st Qu.:0.0000					
	Median : 4467.9	Median :0.0000	Median :0.0000	Median :0.0000					
	Mean : 8468.3	Mean :0.3738	Mean :0.1308	Mean :0.1402					
	3rd Qu.:13627.7	3rd Qu.:1.0000	3rd Qu.:0.0000	3rd Qu.:0.0000					
	Max. :36741.0	Max. :1.0000	Max. :1.0000	Max. :1.0000					
	growth								
	Min. :-0.6888								
	1st Qu.: 0.2458								
	Median : 0.6587								
	Mean : 0.6345								
	3rd Qu.: 1.0505								
	Max. : 2.3493								

```
view(data)
```

```
# library(vtable)
# vtable(data, missing=TRUE)

# library(pastecs)
stat.desc(data)
```

## 9. Collection of exercises

	country	gdppc60	gdppc65	gdppc70	gdppc75	africa
nbr.val	NA	1.070000e+02	1.070000e+02	1.070000e+02	1.070000e+02	
nbr.null	NA	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	
nbr.na	NA	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	
min	NA	4.078180e+02	5.135667e+02	3.545075e+02	6.178639e+02	
max	NA	1.601025e+04	1.892888e+04	2.203095e+04	2.180892e+04	
range	NA	1.560243e+04	1.841531e+04	2.167644e+04	2.119105e+04	
sum	NA	3.888715e+05	4.673224e+05	5.487424e+05	6.162241e+05	
median	NA	2.484720e+03	2.884388e+03	3.072176e+03	3.741725e+03	
mean	NA	3.634313e+03	4.367500e+03	5.128433e+03	5.759103e+03	
SE.mean	NA	3.314566e+02	4.021934e+02	4.736475e+02	5.272377e+02	
CI.mean	NA	6.571449e+02	7.973875e+02	9.390523e+02	1.045300e+03	
var	NA	1.175539e+07	1.730827e+07	2.400459e+07	2.974381e+07	
std.dev	NA	3.428613e+03	4.160321e+03	4.899448e+03	5.453789e+03	
coef.var	NA	9.434006e-01	9.525635e-01	9.553499e-01	9.469857e-01	
	gdppc80	gdppc85	gdppc90	gdppc95		
nbr.val	1.070000e+02	1.070000e+02	1.070000e+02	1.070000e+02	107.000000000	
nbr.null	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	67.000000000	
nbr.na	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000000	
min	4.735793e+02	5.422725e+02	5.277151e+02	4.993415e+02	0.000000000	
max	2.386009e+04	2.525136e+04	2.874414e+04	3.674105e+04	1.000000000	
range	2.338651e+04	2.470909e+04	2.821642e+04	3.624171e+04	1.000000000	
sum	7.012400e+05	7.383373e+05	8.319308e+05	9.061030e+05	40.000000000	
median	4.306217e+03	4.200733e+03	4.034010e+03	4.467940e+03	0.000000000	
mean	6.553645e+03	6.900349e+03	7.775054e+03	8.468253e+03	0.37383178	
SE.mean	6.018749e+02	6.552251e+02	7.711596e+02	8.456513e+02	0.04699273	
CI.mean	1.193276e+03	1.299048e+03	1.528899e+03	1.676586e+03	0.09316766	
var	3.876112e+07	4.593724e+07	6.363152e+07	7.651850e+07	0.23628990	
std.dev	6.225843e+03	6.777701e+03	7.976937e+03	8.747486e+03	0.48609659	
coef.var	9.499817e-01	9.822259e-01	1.025965e+00	1.032974e+00	1.30030838	
	asia	weurope	growth			
nbr.val	107.00000000	107.00000000	107.000000			
nbr.null	93.00000000	92.00000000	0.0000000			
nbr.na	0.00000000	0.00000000	0.0000000			
min	0.00000000	0.00000000	-0.6887722			
max	1.00000000	1.00000000	2.3493433			
range	1.00000000	1.00000000	3.0381155			
sum	14.00000000	15.00000000	67.8899760			
median	0.00000000	0.00000000	0.6586871			
mean	0.13084112	0.14018692	0.6344858			
SE.mean	0.03275433	0.03372119	0.0601857			
CI.mean	0.06493865	0.06685553	0.1193240			
var	0.11479457	0.12167166	0.3875881			
std.dev	0.33881347	0.34881465	0.6225657			
coef.var	2.58950297	2.48821120	0.9812131			

```
# library(Hmisc)
describe(data)
```

vars	n	mean	sd	median	trimmed	mad	min	max
------	---	------	----	--------	---------	-----	-----	-----

## 9. Collection of exercises

country*	1	107	54.00	31.03	54.00	54.00	40.03	1.00	107.00
gdppc60	2	107	3634.31	3428.61	2484.72	3032.19	2027.76	407.82	16010.25
gdppc65	3	107	4367.50	4160.32	2884.39	3673.42	2579.50	513.57	18928.88
gdppc70	4	107	5128.43	4899.45	3072.18	4370.29	2854.11	354.51	22030.95
gdppc75	5	107	5759.10	5453.79	3741.72	4977.54	3708.25	617.86	21808.92
gdppc80	6	107	6553.64	6225.84	4306.22	5707.40	4476.29	473.58	23860.09
gdppc85	7	107	6900.35	6777.70	4200.73	5929.46	4382.44	542.27	25251.36
gdppc90	8	107	7775.05	7976.94	4034.01	6660.00	4258.37	527.72	28744.14
gdppc95	9	107	8468.25	8747.49	4467.94	7235.12	4935.09	499.34	36741.05
africa	10	107	0.37	0.49	0.00	0.34	0.00	0.00	1.00
asia	11	107	0.13	0.34	0.00	0.05	0.00	0.00	1.00
weurope	12	107	0.14	0.35	0.00	0.06	0.00	0.00	1.00
growth	13	107	0.63	0.62	0.66	0.63	0.59	-0.69	2.35
			range	skew	kurtosis	se			
country*			106.00	0.00	-1.23	3.00			
gdppc60			15602.43	1.53	1.55	331.46			
gdppc65			18415.31	1.41	1.16	402.19			
gdppc70			21676.44	1.29	0.74	473.65			
gdppc75			21191.05	1.15	0.09	527.24			
gdppc80			23386.51	1.07	-0.11	601.87			
gdppc85			24709.09	1.14	0.01	655.23			
gdppc90			28216.42	1.13	-0.10	771.16			
gdppc95			36241.71	1.14	0.12	845.65			
africa			1.00	0.51	-1.75	0.05			
asia			1.00	2.16	2.69	0.03			
weurope			1.00	2.04	2.20	0.03			
growth			3.04	0.15	0.07	0.06			

```
# library(gtsummary)
#tbl_summary(data)

# check the assignments of countries to continents
data |>
  select(country, africa, asia, weurope) |>
  view()

data <- mutate(data, x_1 = africa + asia + weurope)

data |>
  filter(x_1 == 0) |>
  select(africa, asia, weurope, country) |>
  view()

# correct the assignment manually
data$weurope[data$country == "Austria"] <- 1
data$weurope[data$country == "Greece"] <- 1
data$weurope[data$country == "Cyprus"] <- 1

filter(data, data$weurope == 1) # check changes
```

## 9. Collection of exercises

```
# A tibble: 18 x 14
  country gdppc60 gdppc65 gdppc70 gdppc75 gdppc80 gdppc85 gdppc90 gdppc95
  <chr>    <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
1 Austria   7842.   9387.  11946.  14198.  16869.  17919.  21178.  22474.
2 Belgium   8314.   10454.  12980.  15024.  17451.  18109.  21246.  22356.
3 Cyprus    3178.   4261.   5638.   4827.   8302.   10228.  13798.  17169.
4 Denmark   11745.  14749.  17143.  17750.  19558.  21596.  23308.  25293.
5 Finland   8007.   9851.   12198.  14884.  16621.  18585.  21667.  20084.
6 France    8364.   10497.  13186.  14951.  17335.  18429.  21403.  21502.
7 Greece    4454.   6549.   9022.   11121.  12672.  12287.  12794.  13332.
8 Iceland   8786.   11403.  11678.  15235.  19440.  20414.  22502.  21901.
9 Ireland   5490.   6413.   7760.   9064.   10649.  11641.  15133.  18456.
10 Italy     7364.   9097.   12072.  13386.  16286.  17518.  20638.  21691.
11 Luxembourg 12510.  14019.  16163.  17384.  19089.  21414.  28744.  36741.
12 Netherlands 9883.  11702.  14237.  15803.  17339.  17974.  20823.  22320.
13 Norway    8808.   10478.  11959.  14873.  17977.  20630.  21855.  25538.
14 Portugal   3665.   4866.   6730.   7951.   9667.   9847.   13155.  13924.
15 Spain      4956.   7459.   9701.   11970.  12294.  12583.  15475.  17434.
16 Sweden     10870.  13552.  15850.  17588.  18348.  20001.  22219.  22122.
17 Switzerland 16010.  18929.  22031.  21809.  23860.  24844.  27931.  26227.
18 United Kingd~ 10341.  11633.  12917.  14072.  15302.  16878.  19585.  20963.
```

# i 5 more variables: africa <dbl>, asia <dbl>, weurope <dbl>, growth <dbl>,  
# x\_1 <dbl>

```
# In the following, I do the same with a loop
# c_europe <- c("Austria", "Greece", "Cyprus")
# sum(data$weurope)                                # check changes
# for (i in c_europe){
#   print(i)
#   data$weurope[data$country == i] <- 1
# }
# sum(data$weurope)                                # check changes
# data$weurope[data$country == "Austria"] # check changes

# create a category for the remaining countries
# use ifelse -- ifelse(condition, result if TRUE, result if FALSE)
data$rest <- ifelse(data$africa == 0 & data$asia == 0 & data$weurope == 0, 1, 0)
data$rest <- set_label(data$rest, label = "1 if not in Africa, W.Europe, or Asia")

# create table with means across country groups
table_gdp <- data |>
  group_by(africa, asia, weurope) |>
  summarise_at(vars(gdppc60:gdppc95), list(name = mean))

data |>
  group_by(africa, asia, weurope) |>
  select(gdppc60:gdppc95) |>
  summarise_all(mean)
```

Adding missing grouping variables: `africa`, `asia`, `weurope`

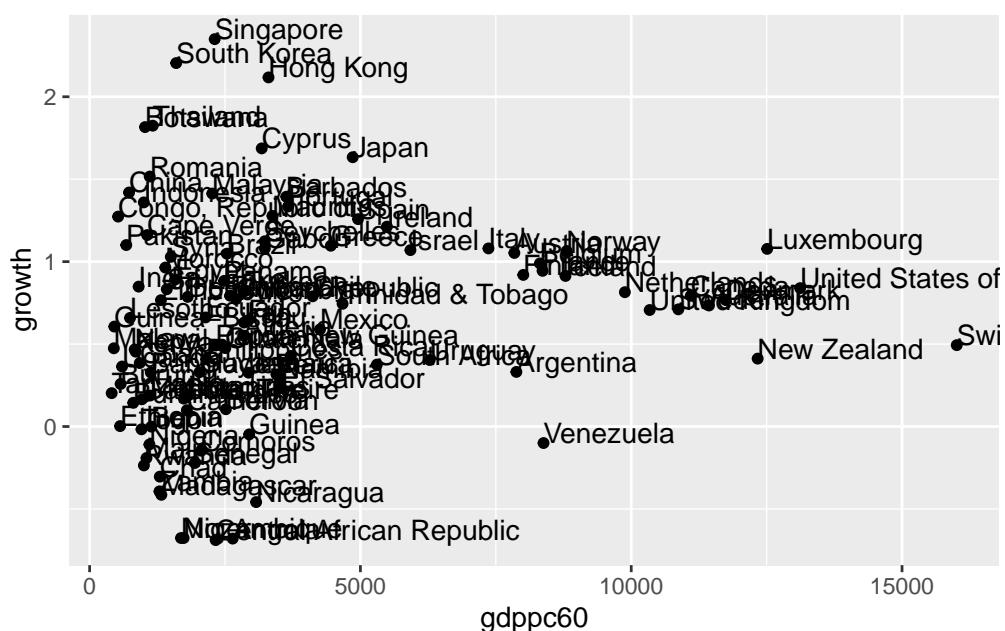
## 9. Collection of exercises

```
# A tibble: 4 x 11
# Groups:   africa, asia [3]
  africa   asia   weurope gdppc60 gdppc65 gdppc70 gdppc75 gdppc80 gdppc85 gdppc90
    <dbl> <dbl>   <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
1      0     0       0    4288.   5034.   5727.   6411.   7042.   7185.   7457.
2      0     0       1   8366.  10294.  12401.  13994.  16059.  17272.  20192.
3      0     1       0   1739.   2247.   3090.   3760.   4905.   5761.   7501.
4      1     0       0   1596.   1860.   2046.   2182.   2426.   2382.   2562.
# i 1 more variable: gdppc95 <dbl>
```

```
# create growth rate
data$gr1 <- (data$gdppc95 - data$gdppc60) / data$gdppc60
data$gr2 <- log(data$gdppc95) - log(data$gdppc60)
cor(data$gr1, data$gr2)
```

[1] 0.9008887

```
ggplot(data, aes(x = gdppc60, y = growth, label = country)) +
  geom_point() +
  geom_text(hjust = 0, vjust = 0)
```



## 9. Collection of exercises

```
p1 <- ggplot(data, aes(x = gdppc60, y = growth, label = country)) +
  geom_point() +
  stat_smooth(formula = y ~ x, method = "lm", se = FALSE, colour = "red", linetype =
  # geom_text(hjust=0, vjust=0) +
  ggtitle("World")

p2 <- data |>
  filter(weurope == 1) |>
  ggplot(aes(x = gdppc60, y = growth, label = country)) +
  geom_point() +
  stat_smooth(formula = y ~ x, method = "lm", se = FALSE, colour = "red", linetype =
  # geom_text(hjust=0, vjust=0) +
  ggtitle("Western Europe")

p3 <- data |>
  filter(asia == 1) |>
  ggplot(aes(x = gdppc60, y = growth, label = country)) +
  geom_point() +
  stat_smooth(formula = y ~ x, method = "lm", se = FALSE, colour = "red", linetype =
  # geom_text(hjust=0, vjust=0) +
  ggtitle("Asia")

p4 <- data |>
  filter(africa == 1) |>
  ggplot(aes(x = gdppc60, y = growth, label = country)) +
  geom_point() +
  stat_smooth(formula = y ~ x, method = "lm", se = FALSE, colour = "red", linetype =
  # geom_text( hjust=0, vjust=0) +
  ggtitle("Africa")

ggarrange(p1, p2, p3, p4,
  labels = c("A", "B", "C", "D"),
  ncol = 2, nrow = 2
)
```

Warning: The following aesthetics were dropped during statistical transformation: label

- i This can happen when ggplot fails to infer the correct grouping structure in the data.
- i Did you forget to specify a `group` aesthetic or to convert a numerical variable into a factor?

The following aesthetics were dropped during statistical transformation: label.

- i This can happen when ggplot fails to infer the correct grouping structure in the data.
- i Did you forget to specify a `group` aesthetic or to convert a numerical variable into a factor?

The following aesthetics were dropped during statistical transformation: label.

- i This can happen when ggplot fails to infer the correct grouping structure in the data.
- i Did you forget to specify a `group` aesthetic or to convert a numerical variable into a factor?

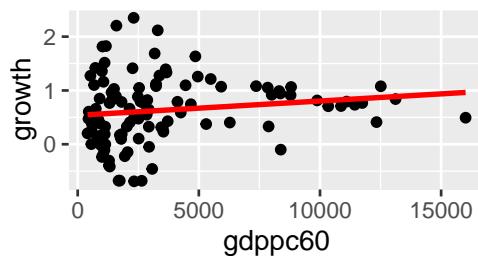
## 9. Collection of exercises

The following aesthetics were dropped during statistical transformation: label.

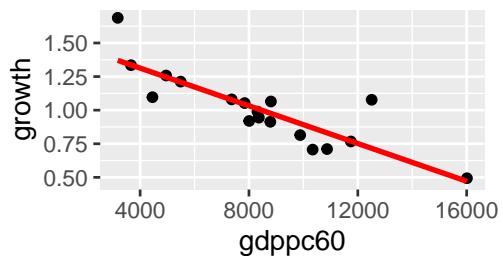
i This can happen when ggplot fails to infer the correct grouping structure in the data.

i Did you forget to specify a `group` aesthetic or to convert a numerical variable into a factor?

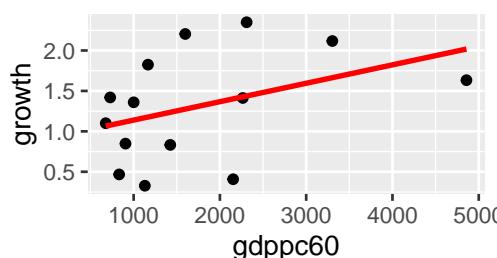
**A World**



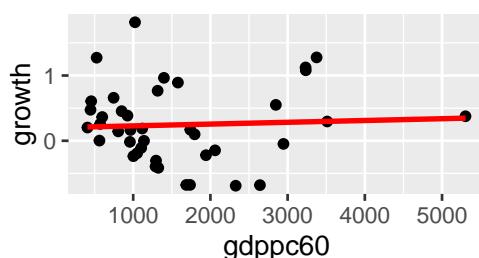
**B Western Europe**



**C Asia**



**D Africa**



```
# Regression analysis
m1 <- lm(growth ~ gdppc60, data = data)
m2 <- lm(growth ~ gdppc60, data = subset(data, weurope == 1))
m3 <- lm(growth ~ gdppc60, data = subset(data, asia == 1))
m4 <- lm(growth ~ gdppc60, data = subset(data, africa == 1))

tab_model(m1, m2, m3, m4,
  p.style = "stars",
  p.threshold = c(0.2, 0.1, 0.05),
  show.ci = FALSE,
  show.se = FALSE,
  show.aic = TRUE,
  dv.labels = c("World", "W.Europe", "Asia", "Africa")
)
```

	World	W.Europe	Asia	Africa
Predictors	Estimates	Estimates	Estimates	Estimates
(Intercept)	0.54 ***	1.59 ***	0.91 ***	0.20
real gdp per capita 1960	0.00 *	-0.00 ***	0.00 *	0.00
Observations	107	18	14	40
R <sup>2</sup> / R <sup>2</sup> adjusted	0.021 / 0.012	0.727 / 0.710	0.158 / 0.088	0.002 / -0.024
AIC	204.917	-14.237	31.220	76.318

\* p<0.2 \*\* p<0.1 \*\*\* p<0.05

## 9. Collection of exercises

```
# reshape data (see: https://stackoverflow.com/questions/2185252/reshaping-data-frame)
data_long <- gather(data, condition, measurement, gdppc60:gdppc95, factor_key = TRUE)

Warning: attributes are not identical across measure variables; they will be
dropped

data_long$year <- as.numeric(substr(data_long$condition, 6, 7))

data_long$gr_long <- data_long |>
  select(country, measurement) |>
  group_by(country) |>
  mutate(gr = c(NA, diff(measurement)) / lag(measurement, 1))

# erase all helping variables
data <- select(data, -starts_with("h_"))

# generate and remove variables in a dataframe
data <- mutate(data, Land = country)
data <- select(data, -country)

data |>
  summarise(
    y65 = mean(gdppc65, na.rm = TRUE),
    y70 = mean(gdppc70, na.rm = TRUE),
    y75 = mean(gdppc75, na.rm = TRUE),
    y80 = mean(gdppc80, na.rm = TRUE),
    y85 = mean(gdppc85, na.rm = TRUE),
    y90 = mean(gdppc90, na.rm = TRUE),
    y95 = mean(gdppc95, na.rm = TRUE)
  )

# A tibble: 1 x 7
#>   y65     y70     y75     y80     y85     y90     y95
#>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 4367. 5128. 5759. 6554. 6900. 7775. 8468.

suppressMessages(pacman::p_unload(
  haven, tidyverse, vtable, gtsummary, pastecs, Hmisc,
  sjlabelled, tis, ggpubr, sjPlot
))
```

## 9.11. Unemployment and GDP in Germany and France

The following exercise was a former exam.

Please answer all (!) questions in an R script. Normal text should be written as comments, using the '#' to comment out text. Make sure the script runs without errors before submitting

## 9. Collection of exercises

it. Each task (starting with 1) is worth five points. You have a total of 120 minutes of editing time. Please do not forget to number your answers.

When you are done with your work, save the R script, export the script to pdf format and upload the pdf file.

Suppose you aim to empirically examine unemployment and GDP for Germany and France. The data set that we use in the following is ‘forest.Rdata’.

- (0) Write down your name, matriculation number, and date.
- (1) Set your working directory.
- (2) Clear your global environment.
- (3) Install and load the following packages: ‘tidyverse’, ‘sjPlot’, and ‘ggpubr’
- (4) Download and load the data, respectively, with the following code:

```
load(url("https://github.com/hubchev/courses/raw/main/dta/forest.Rdata"))
```

If that is not working, you can also download the data from ILIAS, save it in your working directory and load it from there with:

```
# load("forest.Rdata")
```

- (5) Show the **first eight** observations of the dataset **df**.
- (6) Show the **last observation** of the dataset **df**.
- (7) Which type of data do we have here (Panel, cross-section, time series, ...)? Name the variable(s) that are necessary to identify the observations in the dataset.
- (8) Explain what the **assignment operator** in R is and what it is good for.
- (9) Write down the R code to store the number of observations and the number of variables that are in the dataset **df**. Name the object in which you store these numbers ‘**observations\_df**’.
- (10) In the dataset **df**, rename the variable ‘country.x’ to ‘nation’ and the variable ‘date’ to ‘year’.
- (11) Explain what the **pipe operator** in R is and what it is good for.
- (12) For the upcoming analysis you are only interested the following **variables** that are part of the dataframe **df**: nation, year, gdp, pop, gdppc, and unemployment. Drop all other variables from the dataframe **df**.
- (13) Create a variable that indicates the GDP per capita (‘gdp’ divided by ‘pop’). Name the variable ‘**gdp\_pc**’. (Hint: If you fail here, use the variable ‘gdppc’ which is already in the dataset as a replacement for ‘**gdp\_pc**’ in the following tasks.)
- (14) For the upcoming analysis you are only interested the following **countries** that are part of the dataframe **df**: Germany and France. Drop all other countries from the dataframe **df**.
- (15) Create a table showing the **average** unemployment rate and GDP per capita for Germany and France in the given years. Use the pipe operator. (Hint: See below for how your results should look like.)

## 9. Collection of exercises

```
# A tibble: 2 x 3
  nation `mean(unemployment)` `mean(gdppc)`
  <chr>          <dbl>        <dbl>
1 France         9.75       34356.
2 Germany        7.22       36739.
```

- (16) Create a table showing the unemployment rate and GDP per capita for Germany and France in the **year 2020**. Use the pipe operator. (Hint: See below for how your results should look like.)

```
# A tibble: 2 x 3
  nation `mean(unemployment)` `mean(gdppc)`
  <chr>          <dbl>        <dbl>
1 France         8.01       35786.
2 Germany        3.81       41315.
```

- (17) Create a table showing the **highest** unemployment rate and the **highest** GDP per capita for Germany and France during the given period. Use the pipe operator. (Hint: See below for how your results should look like.)

```
# A tibble: 2 x 3
  nation `max(unemployment)` `max(gdppc)`
  <chr>          <dbl>        <dbl>
1 France         12.6       38912.
2 Germany        11.2       43329.
```

- (18) Calculate the standard deviation of the unemployment rate and GDP per capita for Germany and France in the given years. (Hint: See below for how your result should look like.)

```
# A tibble: 2 x 3
  nation `sd(gdppc)` `sd(unemployment)`
  <chr>      <dbl>        <dbl>
1 France     2940.        1.58
2 Germany    4015.        2.37
```

- (19) In statistics, the coefficient of variation (COV) is a standardized measure of dispersion. It is defined as the ratio of the standard deviation ( $\sigma$ ) to the mean ( $\mu$ ):  $COV = \frac{\sigma}{\mu}$ . Write down the R code to calculate the coefficient of variation (COV) for the **unemployment rate** in Germany and France. (Hint: See below for what your result should look like.)

```
# A tibble: 2 x 4
  nation `sd(unemployment)` `mean(unemployment)` cov
  <chr>      <dbl>        <dbl> <dbl> <dbl>
1 France     1.58           9.75  0.162
2 Germany    2.37           7.22  0.328
```

- (20) Write down the R code to calculate the coefficient of variation (COV) for the **GDP per capita** in Germany and France. (Hint: See below for what your result should look like.)

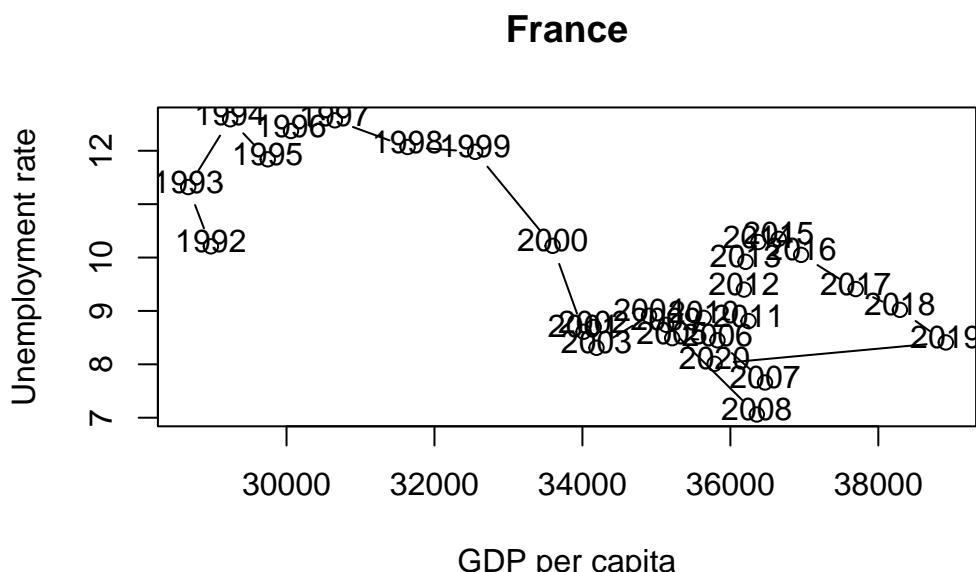
## 9. Collection of exercises

```
# A tibble: 2 x 4
  nation `sd(gdppc)` `mean(gdppc)` cov
  <chr>     <dbl>        <dbl>   <dbl>
1 France      2940.       34356.  0.0856
2 Germany     4015.       36739.  0.109
```

- (21) Create a chart (bar chart, line chart, or scatter plot) that shows the unemployment rate of **Germany** over the available years. Label the chart ‘Germany’ with ‘`ggtitle("Germany")`’. Please note that you may choose any type of graphical representation. (Hint: Below you can see one of many `|>` of what your result may look like).



- (22) and 23. (*This task is worth 10 points*) The following chart shows the simultaneous development of the unemployment rate and GDP per capita over time for **France**.



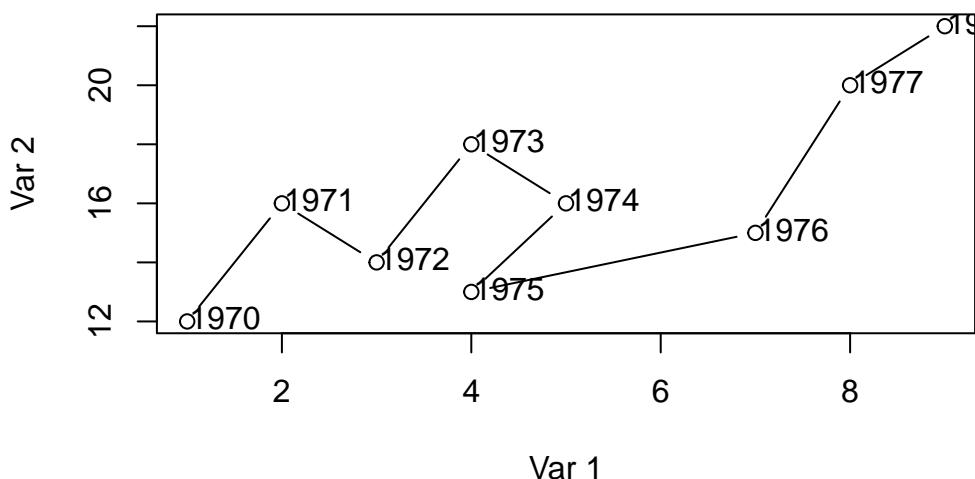
Suppose you want to visualize the simultaneous evolution of the unemployment rate and GDP per capita over time for Germany as well.

## 9. Collection of exercises

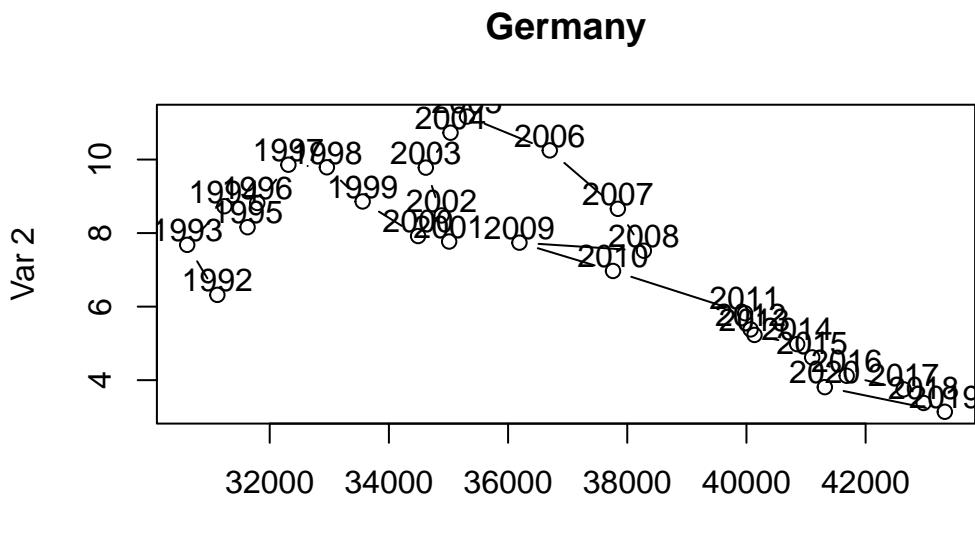
Suppose further that you have found the following lines of code that create the kind of chart you are looking for.

```
# Data
x <- c(1, 2, 3, 4, 5, 4, 7, 8, 9)
y <- c(12, 16, 14, 18, 16, 13, 15, 20, 22)
labels <- 1970:1978

# Connected scatter plot with text
plot(x, y, type = "b", xlab = "Var 1", ylab = "Var 2")
text(x + 0.4, y + 0.1, labels)
```



Use these lines of code and customize them to create the co-movement visualization for **Germany** using the available `df` data. The result should look something like this:



- (24) Interpret the two graphs above, which show the simultaneous evolution of the unemployment rate and GDP per capita over time for Germany and France. What are your expectations regarding the correlation between the unemployment rate and GDP per capita variables? Can you see this expectation in the figures? Discuss.

 Solution

The script uses the following functions: `aes`, `c`, `dim`, `filter`, `geom_line`, `ggplot`, `ggttitle`, `group_by`, `head`, `load`, `max`, `mean`, `mutate`, `plot`, `rename`, `sd`, `select`, `summarise`, `tail`, `text`, `title`, `url`.

**R script**

## 9. Collection of exercises

```
# setwd("/home/sthu/Dropbox/hsf/exams/22-11/scr/")

rm(list = ls())

if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, ggpunr, sjPlot)

load(url("https://github.com/hubchev/courses/raw/main/dta/forest.Rdata"))

head(df, 8)

tail(df, 1)

# panel data set
# date and country.x

observations_df <- dim(df)

df <- rename(df, nation = country.x)
df <- rename(df, year = date)

df <- df |>
  select(nation, year, gdp, pop, gdppc, unemployment)

df <- df |>
  mutate(gdp_pc = gdp / pop)

df <- df |>
  filter(nation == "Germany" | nation == "France")

df |>
  group_by(nation) |>
  summarise(mean(unemployment), mean(gdppc))

df |>
  filter(year == 2020) |>
  group_by(nation) |>
  summarise(mean(unemployment), mean(gdppc))

df |>
  group_by(nation) |>
  summarise(max(unemployment), max(gdppc))

df |>
  group_by(nation) |>
  summarise(sd(gdppc), sd(unemployment))

df |>
  group_by(nation) |>
  summarise(sd(unemployment), mean(unemployment), cov = sd(unemployment) / mean(unemployment))

df |>
  group_by(nation) |>
  summarise(sd(gdppc), mean(gdppc), cov = sd(gdppc) / mean(gdppc))

df |>
  filter(nation == "Germany") |>
```

## Output of the R script

```
# setwd("/home/sthu/Dropbox/hsf/exams/22-11/scr/")

rm(list = ls())

if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, ggpunr, sjPlot)

load(url("https://github.com/hubchev/courses/raw/main/dta/forest.Rdata"))

head(df, 8)

# A tibble: 8 x 11
# Groups:   country.x [1]
  country.x     date     gdp gdp_growth unemployment region income forest    pop
  <chr>        <dbl>   <dbl>      <dbl>       <dbl> <chr> <chr>   <dbl>   <dbl>
1 United Arab~ 1992 1.26e11    -2.48       1.84 Middl~ High ~  3.63 2.05e6
2 United Arab~ 1993 1.27e11    -4.34       1.85 Middl~ High ~  3.72 2.17e6
3 United Arab~ 1994 1.36e11     1.25       1.81 Middl~ High ~  3.81 2.29e6
4 United Arab~ 1995 1.45e11     1.35       1.80 Middl~ High ~  3.90 2.42e6
5 United Arab~ 1996 1.54e11     0.631      1.90 Middl~ High ~  3.99 2.54e6
6 United Arab~ 1997 1.66e11     2.83       1.98 Middl~ High ~  4.08 2.67e6
7 United Arab~ 1998 1.67e11    -4.77       2.14 Middl~ High ~  4.18 2.81e6
8 United Arab~ 1999 1.72e11    -2.40       2.22 Middl~ High ~  4.27 2.97e6
# i 2 more variables: unemployment_dif <dbl>, gdppc <dbl>

tail(df, 1)

# A tibble: 1 x 11
# Groups:   country.x [1]
  country.x     date     gdp gdp_growth unemployment region income forest    pop
  <chr>        <dbl>   <dbl>      <dbl>       <dbl> <chr> <chr>   <dbl>   <dbl>
1 Zimbabwe    2020 1.94e10    -7.62       5.35 Sub-S~ Lower~  45.1 1.49e7
# i 2 more variables: unemployment_dif <dbl>, gdppc <dbl>
```

## 9. Collection of exercises

```
# panel data set
# date and country.x

observations_df <- dim(df)

df <- rename(df, nation = country.x)
df <- rename(df, year = date)

df <- df |>
  select(nation, year, gdp, pop, gdppc, unemployment)

df <- df |>
  mutate(gdp_pc = gdp / pop)

df <- df |> filter(nation == "Germany" | nation == "France")

df |>
  group_by(nation) |>
  summarise(mean(unemployment), mean(gdppc))

# A tibble: 2 x 3
  nation `mean(unemployment)` `mean(gdppc)`
  <chr>          <dbl>        <dbl>
1 France           9.75       34356.
2 Germany          7.22       36739.

df |>
  filter(year == 2020) |>
  group_by(nation) |>
  summarise(mean(unemployment), mean(gdppc))

# A tibble: 2 x 3
  nation `mean(unemployment)` `mean(gdppc)`
  <chr>          <dbl>        <dbl>
1 France           8.01       35786.
2 Germany          3.81       41315.

df |>
  group_by(nation) |>
  summarise(max(unemployment), max(gdppc))

# A tibble: 2 x 3
  nation `max(unemployment)` `max(gdppc)`
  <chr>          <dbl>        <dbl>
1 France           12.6       38912.
2 Germany          11.2       43329.

df |>
  group_by(nation) |>
  summarise(sd(gdppc), sd(unemployment))
```

## 9. Collection of exercises

```
# A tibble: 2 x 3
  nation `sd(gdppc)` `sd(unemployment)`
  <chr>      <dbl>           <dbl>
1 France       2940.          1.58
2 Germany      4015.          2.37

df |>
  group_by(nation) |>
  summarise(sd(unemployment), mean(unemployment), cov = sd(unemployment) / mean(unemployment))

# A tibble: 2 x 4
  nation `sd(unemployment)` `mean(unemployment)` cov
  <chr>      <dbl>           <dbl> <dbl>
1 France       1.58            9.75 0.162
2 Germany      2.37            7.22 0.328

df |>
  group_by(nation) |>
  summarise(sd(gdppc), mean(gdppc), cov = sd(gdppc) / mean(gdppc))

# A tibble: 2 x 4
  nation `sd(gdppc)` `mean(gdppc)` cov
  <chr>      <dbl>           <dbl> <dbl>
1 France       2940.          34356. 0.0856
2 Germany      4015.          36739. 0.109

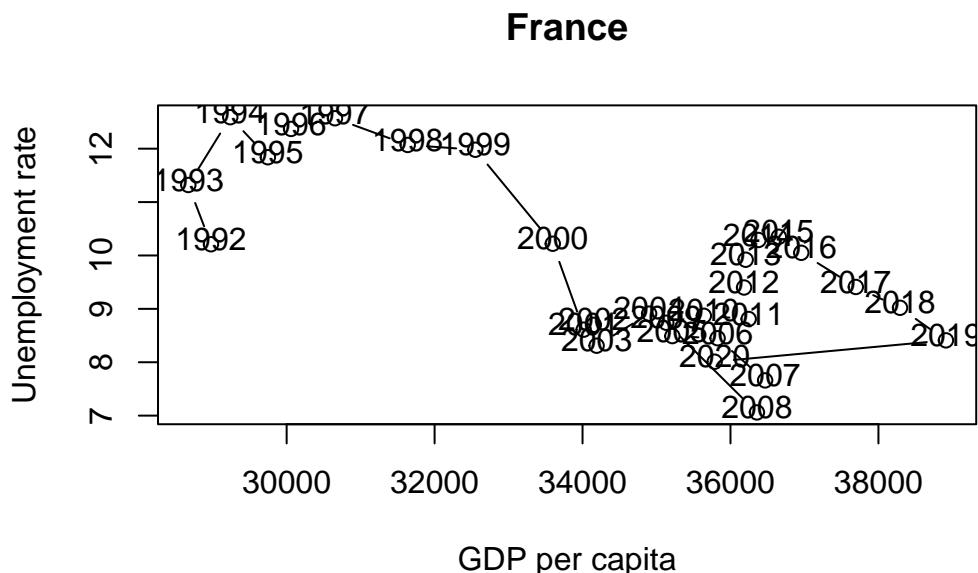
df |>
  filter(nation == "Germany") |>
  ggplot(aes(x = year, y = unemployment)) +
  geom_line() +
  ggtitle("Germany")
```



```

labels <- 1992:2020
dfra <- df |> filter(nation == "France")
plot(dfra$gdppc, dfra$unemployment,
    type = "b",
    xlab = "GDP per capita", ylab = "Unemployment rate"
)
text(dfra$gdppc + 0.1, dfra$unemployment + 0.1, labels)
title("France")

```

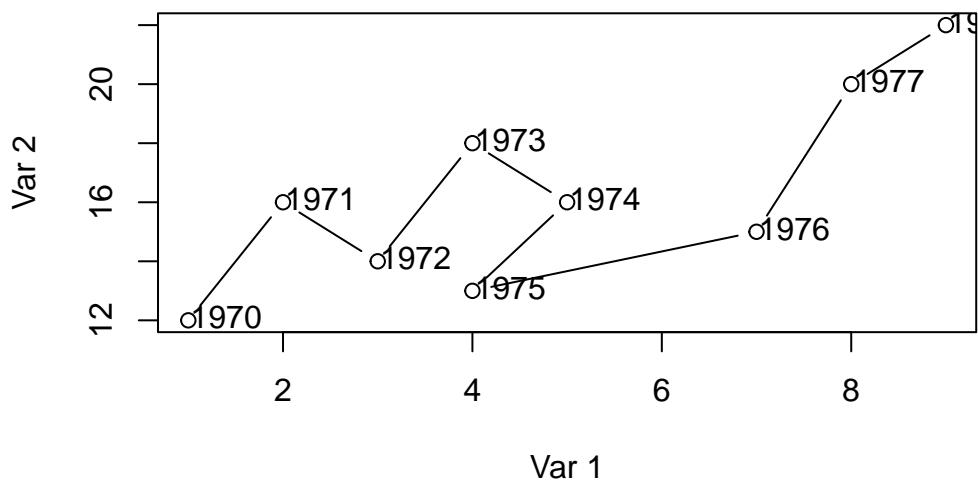


```

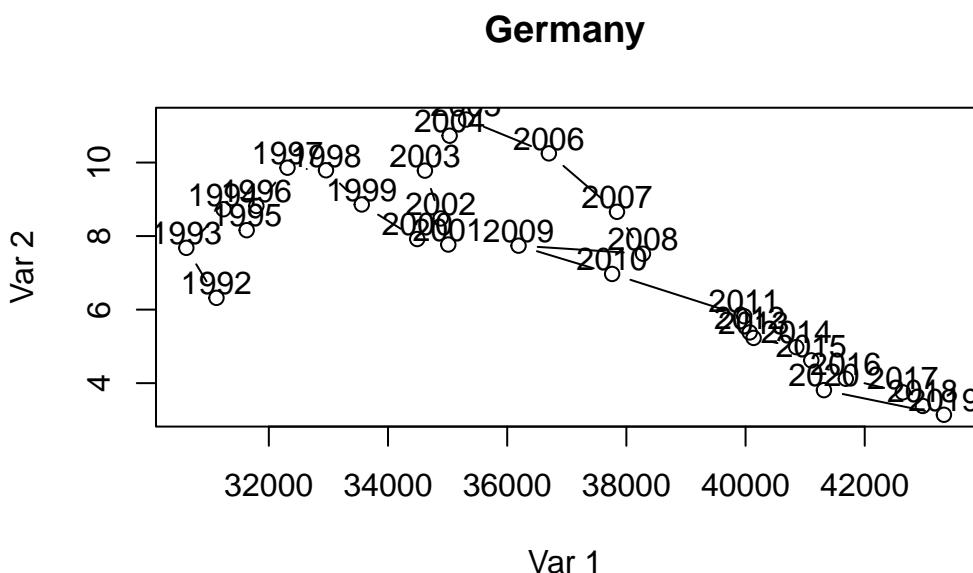
# Data
x <- c(1, 2, 3, 4, 5, 4, 7, 8, 9)
y <- c(12, 16, 14, 18, 16, 13, 15, 20, 22)
labels <- 1970:1978

# Connected scatter plot with text
plot(x, y, type = "b", xlab = "Var 1", ylab = "Var 2")
text(x + 0.4, y + 0.1, labels)

```



```
dfger <- df |> filter(nation == "Germany")
labels <- 1992:2020
plot(dfger$gdppc, dfger$unemployment,
      type = "b",
      xlab = "Var 1", ylab = "Var 2"
)
text(dfger$gdppc + 0.7, dfger$unemployment + 0.4, labels)
title("Germany")
```



```
# rmarkdown::render("22-11_dsda_exam.Rmd", "all")

# knitr::purl(input = "22-11_dsda_exam.Rmd", output = "22-11_dsda_solution.R", document)

suppressMessages(pacman::p_unload(tidyverse, ggpublisher, sjPlot))
```

## 9.12. Import data and write a report

Reproduce Figure 3 of [Hortaçsu and Syverson \[2015, p. 99\]](#) using R. Write a clear report about your work, i.e., document everything with a R script or a R Markdown file.

Here are the required steps:

1. Go to <https://www.aeaweb.org/articles?id=10.1257/jep.29.4.89> and download the *replication package* from the OPENICPSR page. Please note, that you can download the replication package after you have registered for the platform.
2. Unzip the replication package.
3. In the file *diffusion\_curves\_figure.xlsx* you find the required data. Import them to R.
4. Reproduce the plot using `ggplot()`.

 Solution

The script uses the following functions: `aes`, `download.file`, `geom_line`, `ggplot`, `pivot_longer`, `read_excel`, `unzip`.

**R script**

```
# setwd("~/Dropbox/hsf/courses/Rlang/hortacsu")

rm(list = ls())

# install and load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, readxl)

# Define the URL of the ZIP file
zip_f <- "https://github.com/hubchev/courses/raw/main/dta/113962-V1.zip"

# Download the ZIP file
download.file(zip_f, destfile = "113962-V1.zip")

# Unzip the contents
unzip("113962-V1.zip")

df_curves <- read_excel("Hortacsu_Syverson_JEP_Retail/diffusion_curves_figure.xlsx",
  sheet = "Data and Predictions", range = "N3:Y60"
)

df <- df_curves |>
  pivot_longer(
    cols = "Music and Video": "Food and Beverages",
    names_to = "industry",
    values_to = "value"
  )

# Plot
df |>
  ggplot(aes(x = Year, y = value, group = industry, color = industry)) +
  geom_line()

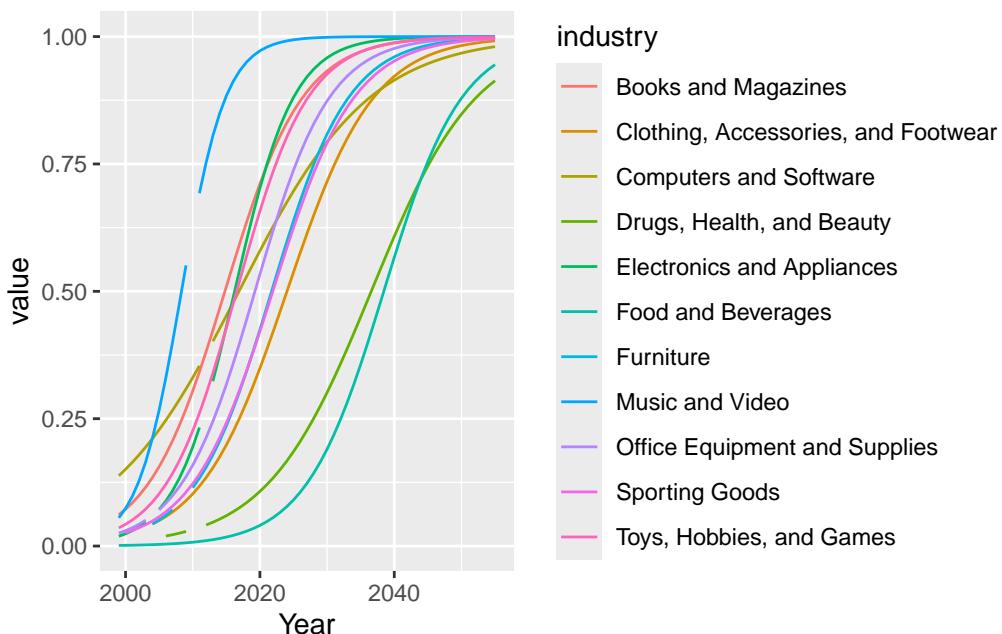
# unload packages
suppressMessages(pacman::p_unload(tidyverse, readxl))
```

**Output of the R script**

## 9. Collection of exercises

```
# setwd("~/Dropbox/hsf/courses/Rlang/hortacsu")  
  
rm(list = ls())  
  
# install and load packages  
if (!require(pacman)) install.packages("pacman")  
pacman::p_load(tidyverse, readxl)  
  
# Define the URL of the ZIP file  
zip_f <- "https://github.com/hubchev/courses/raw/main/dta/113962-V1.zip"  
  
# Download the ZIP file  
download.file(zip_f, destfile = "113962-V1.zip")  
  
# Unzip the contents  
unzip("113962-V1.zip")  
  
df_curves <- read_excel("Hortacsu_Syverson_JEP_Retail/diffusion_curves_figure.xlsx",  
  sheet = "Data and Predictions", range = "N3:Y60"  
)  
  
df <- df_curves |>  
  pivot_longer(  
    cols = "Music and Video": "Food and Beverages",  
    names_to = "industry",  
    values_to = "value"  
  )  
  
# Plot  
df |>  
  ggplot(aes(x = Year, y = value, group = industry, color = industry)) +  
  geom_line()
```

Warning: Removed 18 rows containing missing values or values outside the scale range (`geom\_line()`).



```
# unload packages
suppressMessages(pacman::p_unload(tidyverse, readxl))
```

### 9.13. Explain the weight of students

In the statistic course of WS 2020, I asked 23 students about their weight, height, sex, and number of siblings. I wonder how good the height can explain the weight of students. Examine with corelations and a regression analysis the association. Load the data as follows:

```
library("haven")
```

Attaching package: 'haven'

The following objects are masked from 'package:expss':

```
is.labelled, read_spss
```

```
classdata <- read.csv("https://raw.githubusercontent.com/hubchev/courses/main/dta/classdata
```

#### Solution

The script uses the following functions: `aes`, `c`, `coef`, `fitted`, `geom_abline`, `geom_point`, `ggplot`, `head`, `lm`, `plot`, `read.csv`, `residuals`, `show`, `stat_smooth`, `subset`, `summary`, `tab_model`.

#### R script

## 9. Collection of exercises

```
## ---- echo = TRUE-----
# install and load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, haven, ggplot2, sjPlot)

classdata <- read.csv("https://raw.githubusercontent.com/hubchev/courses/main/dta/cla

head(classdata)

## ---- echo = TRUE-----

summary(classdata)

## ----pressure, echo=TRUE-----
ggplot(classdata, aes(x = height, y = weight)) +
  geom_point()

## ---- echo=TRUE-----
ggplot(classdata, aes(x = height, y = weight)) +
  geom_point() +
  stat_smooth(formula = y ~ x, method = "lm", se = FALSE, colour = "red", linetype = 1)

## ---- echo=TRUE-----
## baseline regression model
model <- lm(weight ~ height + sex, data = classdata)
show(model)
interm <- model$coefficients[1]
slope <- model$coefficients[2]
interw <- model$coefficients[1] + model$coefficients[3]

## ---- echo=TRUE-----
summary(model)

## ---- echo=TRUE-----
ggplot(classdata, aes(x = height, y = weight, shape = sex)) +
  geom_point() +
  geom_abline(slope = slope, intercept = interw, linetype = 2, size = 1.5) +
  geom_abline(slope = slope, intercept = interm, linetype = 2, size = 1.5) +
  geom_abline(slope = coef(model)[[2]], intercept = coef(model)[[1]])

## ---- echo=TRUE-----
ggplot(classdata, aes(x = height, y = weight, shape = sex)) +
  geom_point(aes(size = 2)) +
  stat_smooth(
    formula = y ~ x, method = "lm",
    se = FALSE, colour = "red", linetype = 1
  )

## ---- echo=TRUE-----
ggplot(classdata, aes(x = height177, y = weight, shape = sex)) +
  geom_point(aes(size = siblings))
```

### Output of the R script

```
## ---- echo = TRUE-----
# install and load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, haven, ggplot2, sjPlot)

classdata <- read.csv("https://raw.githubusercontent.com/hubchev/courses/main/dta/classdata.csv")

head(classdata)

   id sex weight height siblings row
1  1   w     53     156       1   g
2  2   w     73     170       1   g
3  3   m     68     169       1   g
4  4   w     67     166       1   g
5  5   w     65     175       1   g
6  6   w     48     161       0   g

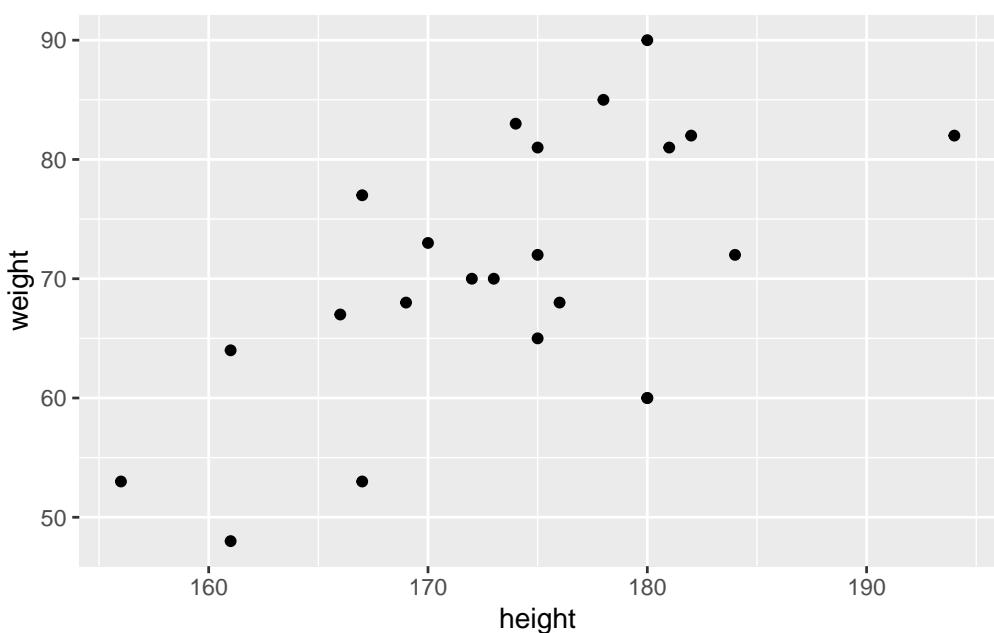
## ---- echo = TRUE-----
summary(classdata)

      id          sex         weight        height
Min. : 1.0  Length:23  Min.   :48.00  Min.   :156.0
1st Qu.: 6.5  Class  :character  1st Qu.:64.50  1st Qu.:168.0
Median :12.0  Mode   :character  Median :70.00  Median :175.0
Mean   :12.0                           Mean   :70.61  Mean   :173.7
3rd Qu.:17.5                           3rd Qu.:81.00  3rd Qu.:180.0
Max.   :23.0                           Max.   :90.00  Max.   :194.0

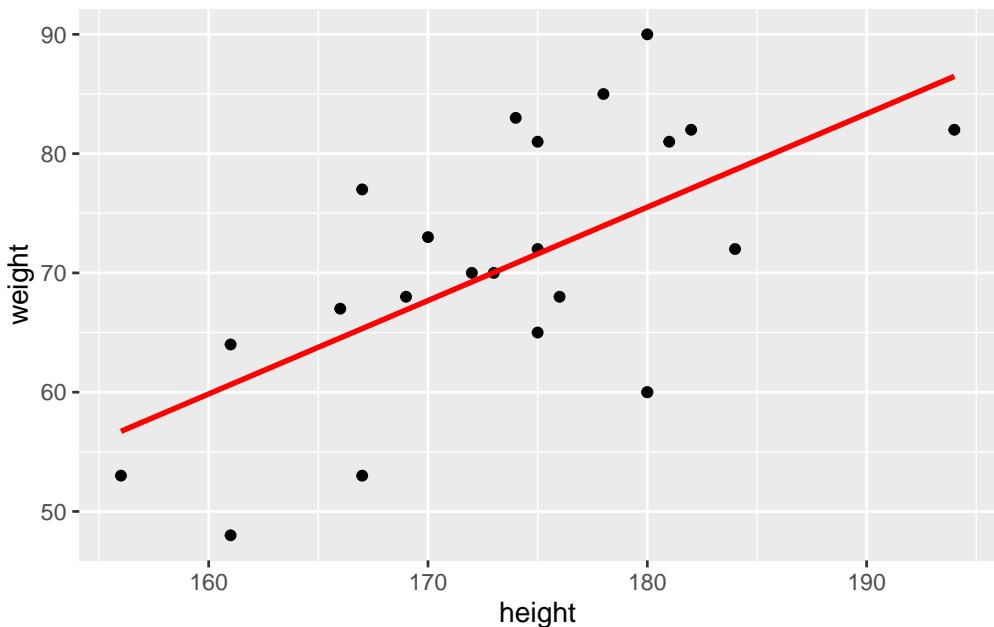
      siblings        row
Min.   :0.000  Length:23
1st Qu.:1.000  Class  :character
Median :1.000  Mode   :character
Mean   :1.391
3rd Qu.:2.000
Max.   :4.000

## ----pressure, echo=TRUE-----
ggplot(classdata, aes(x = height, y = weight)) +
  geom_point()
```

## 9. Collection of exercises



```
## ----- echo=TRUE -----
ggplot(classdata, aes(x = height, y = weight)) +
  geom_point() +
  stat_smooth(formula = y ~ x, method = "lm", se = FALSE, colour = "red", linetype =
```



```
## ----- echo=TRUE -----
## baseline regression model
model <- lm(weight ~ height + sex, data = classdata)
show(model)
```

Call:  
`lm(formula = weight ~ height + sex, data = classdata)`

```

Coefficients:
(Intercept)      height       sexw
-29.5297        0.5923      -5.7894

interm <- model$coefficients[1]
slope <- model$coefficients[2]
interw <- model$coefficients[1] + model$coefficients[3]

## ---- echo=TRUE-----
summary(model)

Call:
lm(formula = weight ~ height + sex, data = classdata)

Residuals:
    Min      1Q  Median      3Q     Max 
-17.086 -3.730   2.850   7.245  12.914 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -29.5297    47.6606  -0.620   0.5425    
height        0.5923     0.2671   2.217   0.0383 *  
sexw         -5.7894    4.4773  -1.293   0.2107    
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

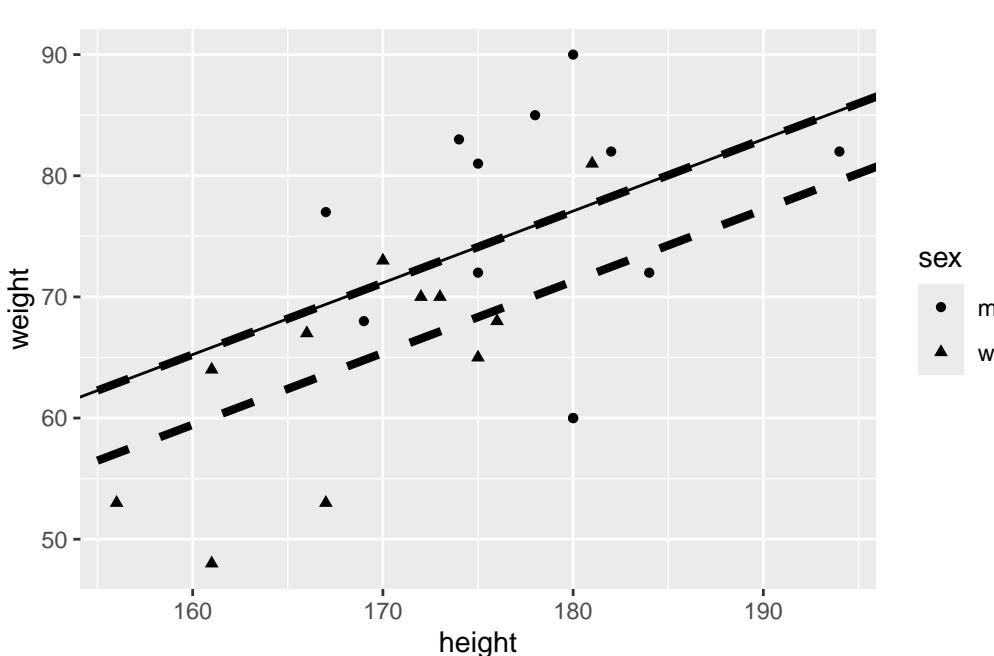
Residual standard error: 8.942 on 20 degrees of freedom
Multiple R-squared:  0.4124,    Adjusted R-squared:  0.3537 
F-statistic: 7.019 on 2 and 20 DF,  p-value: 0.004904

## ---- echo=TRUE-----
ggplot(classdata, aes(x = height, y = weight, shape = sex)) +
  geom_point() +
  geom_abline(slope = slope, intercept = interw, linetype = 2, size = 1.5) +
  geom_abline(slope = slope, intercept = interm, linetype = 2, size = 1.5) +
  geom_abline(slope = coef(model)[[2]], intercept = coef(model)[[1]])

Warning: Using `size` aesthetic for lines was deprecated in ggplot2 3.4.0.
  i Please use `linewidth` instead.

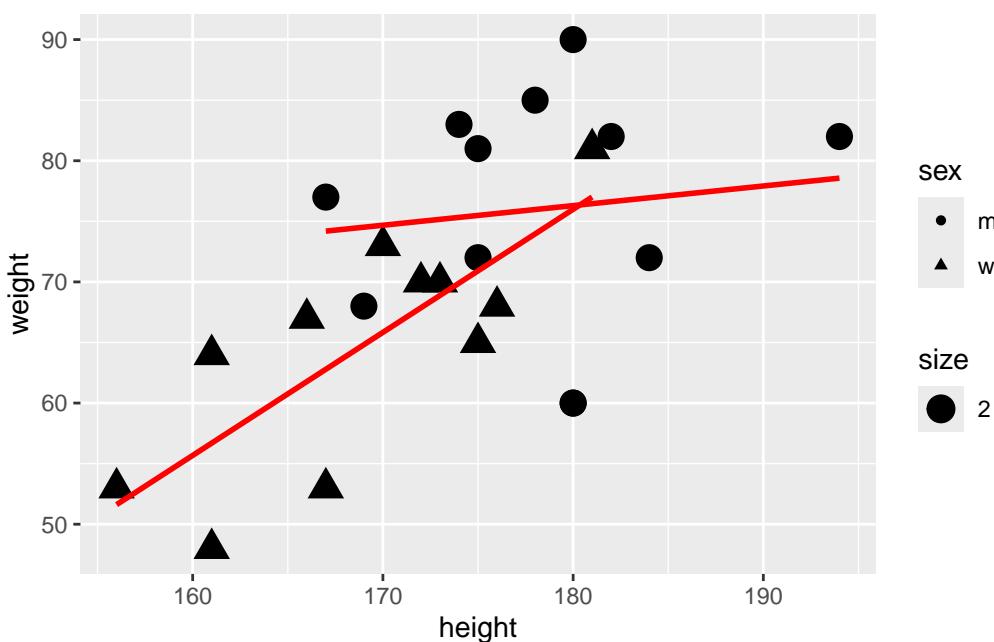
```

## 9. Collection of exercises



```
## ----- echo=TRUE -----

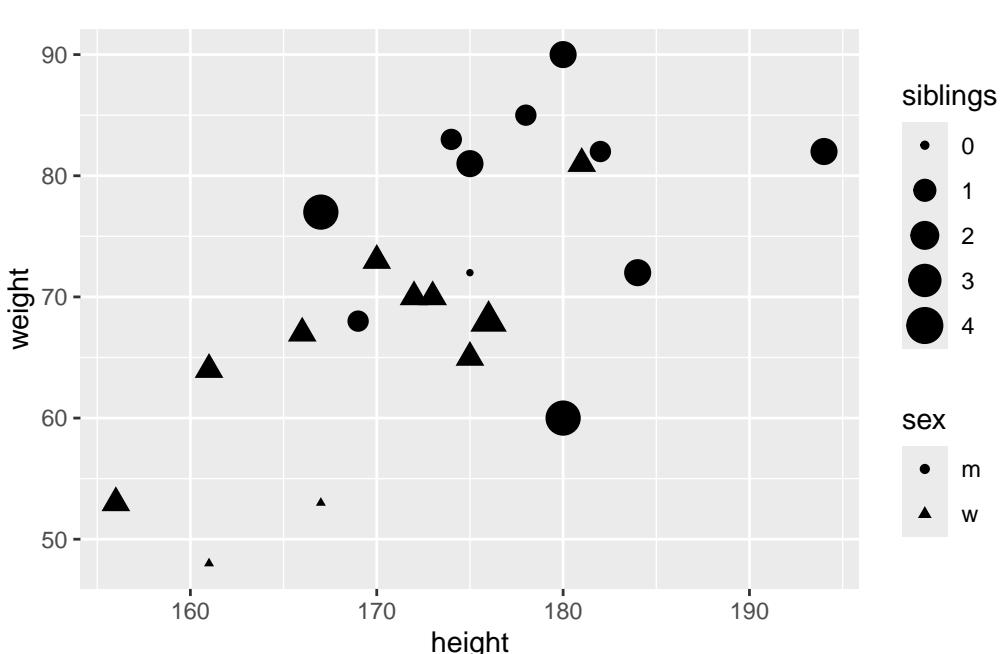
ggplot(classdata, aes(x = height, y = weight, shape = sex)) +
  geom_point(aes(size = 2)) +
  stat_smooth(
    formula = y ~ x, method = "lm",
    se = FALSE, colour = "red", linetype = 1
  )
```



```
## ----- echo=TRUE -----

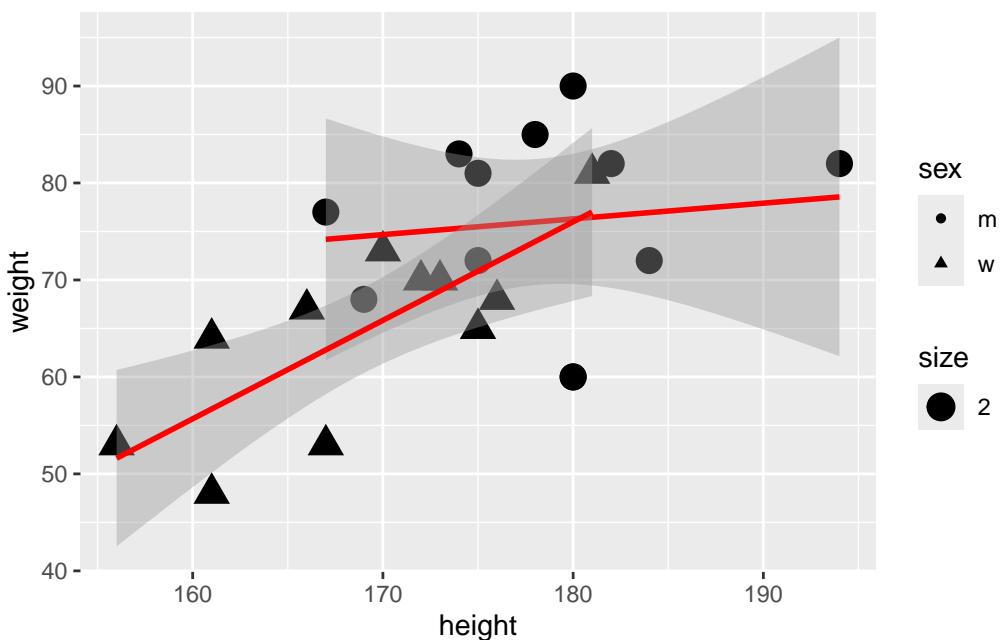
ggplot(classdata, aes(x = height, y = weight, shape = sex)) +
  geom_point(aes(size = siblings))
```

## 9. Collection of exercises



```
## ----- echo=TRUE -----
## baseline model
model <- lm(weight ~ height + sex, data = classdata)

ggplot(classdata, aes(x = height, y = weight, shape = sex)) +
  geom_point(aes(size = 2)) +
  stat_smooth(
    formula = y ~ x,
    method = "lm",
    se = T,
    colour = "red",
    linetype = 1
  )
```



## 9. Collection of exercises

```
## ---- echo=TRUE, results='hide'-----

m1 <- lm(weight ~ height, data = classdata)
m2 <- lm(weight ~ height + sex, data = classdata)
m3 <- lm(weight ~ height + sex + height * sex, data = classdata)
m4 <- lm(weight ~ height + sex + height * sex + siblings, data = classdata)
m5 <- lm(weight ~ height + sex + height * sex, data = subset(classdata, siblings < 4))

tab_model(m1, m2, m3, m4, m5,
  p.style = "stars",
  p.threshold = c(0.2, 0.1, 0.05),
  show.ci = FALSE,
  show.se = FALSE
)
```

Predictors	weight Estimates	weight Estimates	weight Estimates	weight Estimates	weight Estimates
(Intercept)	-65.44 *	-29.53	47.14	50.27	27.69
height	0.78 ***	0.59 ***	0.16	0.16	0.28
sex [w]		-5.79	-153.96 **	-161.92 **	-134.51 *
height × sex [w]			0.85 *	0.89 *	0.74 *
siblings				-1.16	
Observations	23	23	23	23	21
R <sup>2</sup> / R <sup>2</sup> adjusted	0.363 / 0.333	0.412 / 0.354	0.487 / 0.407	0.496 / 0.385	0.572 / 0.497

\* p<0.2    \*\* p<0.1    \*\*\* p<0.05

## 9. Collection of exercises

```
## ----- echo=FALSE -----
tab_model(m1, m2, m3, m4,
  p.style = "stars",
  p.threshold = c(0.2, 0.1, 0.05),
  show.ci = FALSE,
  show.se = FALSE
)
```

Predictors	weight Estimates	weight Estimates	weight Estimates	weight Estimates
(Intercept)	-65.44 *	-29.53	47.14	50.27
height	0.78 ***	0.59 ***	0.16	0.16
sex [w]		-5.79	-153.96 **	-161.92 **
height × sex [w]			0.85 *	0.89 *
siblings				-1.16
Observations	23	23	23	23
R <sup>2</sup> / R <sup>2</sup> adjusted	0.363 / 0.333	0.412 / 0.354	0.487 / 0.407	0.496 / 0.385

\* p<0.2    \*\* p<0.1    \*\*\* p<0.05

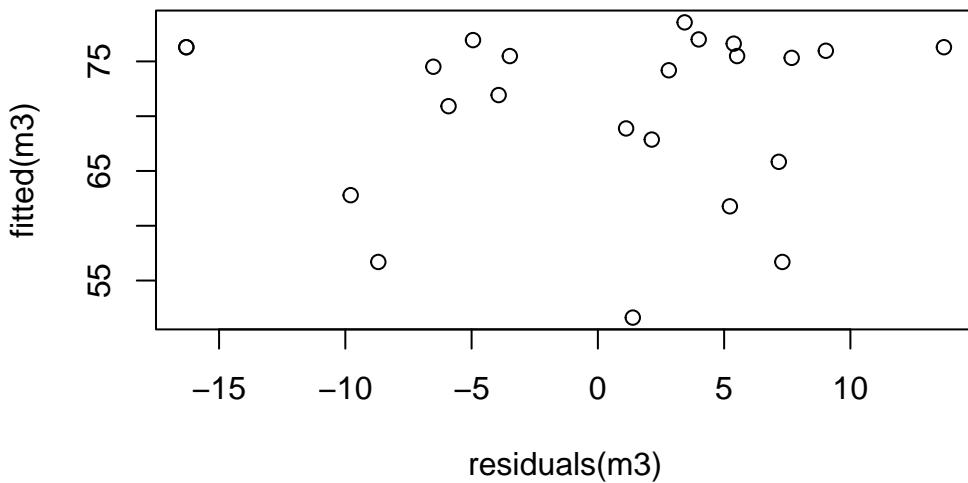
## 9. Collection of exercises

```
## ----- echo=FALSE -----
tab_model(m3, m5,
  p.style = "stars",
  p.threshold = c(0.2, 0.1, 0.05),
  show.ci = FALSE,
  show.se = FALSE
)
```

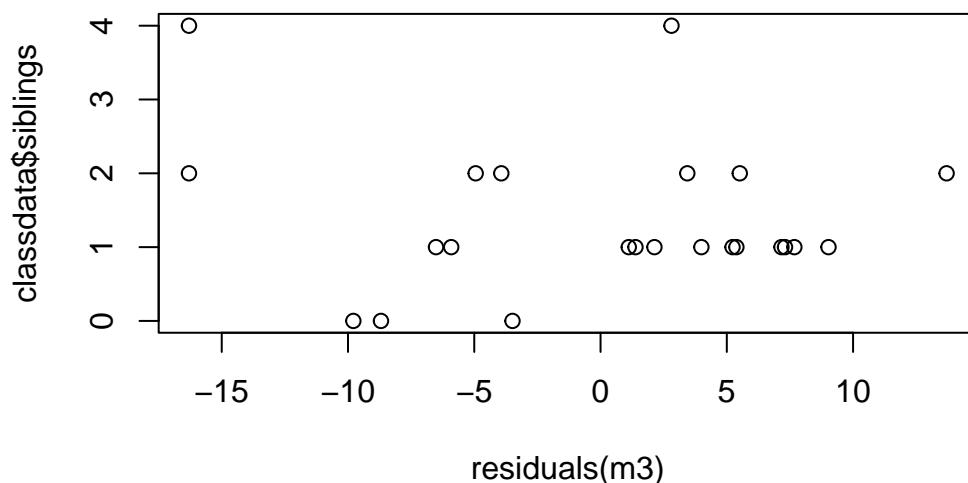
Predictors	Estimates	Estimates
(Intercept)	47.14	27.69
height	0.16	0.28
sex [w]	-153.96 **	-134.51 *
height × sex [w]	0.85 *	0.74 *
Observations	23	21
R <sup>2</sup> / R <sup>2</sup> adjusted	0.487 / 0.407	0.572 / 0.497

\* p<0.2   \*\* p<0.1   \*\*\* p<0.05

```
## ----- echo=T -----
plot(residuals(m3), fitted(m3))
```



```
plot(residuals(m3), classdata$siblings)
```



```
# unload packages
suppressMessages(pacman::p_unload(tidyverse, haven, sjPlot))
```

## 9.14. Calories and weight

- Write down your name, your matriculation number, and the date.
- Set your working directory.
- Clear your global environment.
- Load the following package: `tidyverse`

The following table stems from a survey carried out at the Campus of the German Sport University of Cologne at Opening Day (first day of the new semester) between 8:00am and 8:20am. The survey consists of 6 individuals with the following information:

## 9. Collection of exercises

id	sex	age	weight	calories	sport
1	f	21	48	1700	60
2	f	19	55	1800	120
3	f	23	50	2300	180
4	m	18	71	2000	60
5	m	20	77	2800	240
6	m	61	85	2500	30

Data Description:

- **id:** Variable with an anonymized identifier for each participant.
  - **sex:** Gender, i.e., the participants replied to be either male (m) or female (f).
  - **age:** The age in years of the participants at the time of the survey.
  - **weight:** Number of kg the participants pretended to weight.
  - **calories:** Estimate of the participants on their average daily consumption of calories.
  - **sport:** Estimate of the participants on their average daily time that they spend on doing sports (measured in minutes).
- e) Which type of data do we have here? (Panel data, repeated cross-sectional data, cross-sectional data, time Series data)
- f) Store each of the five variables in a vector and put all five variables into a dataframe with the title `df`. If you fail here, read in the data using this line of code:

```
df <- read_csv("https://raw.githubusercontent.com/hubchev/courses/main/dta/df-calories.csv")
```

Rows: 6 Columns: 5

-- Column specification -----

Delimiter: ","

chr (1): sex

dbl (4): age, weight, calories, sport

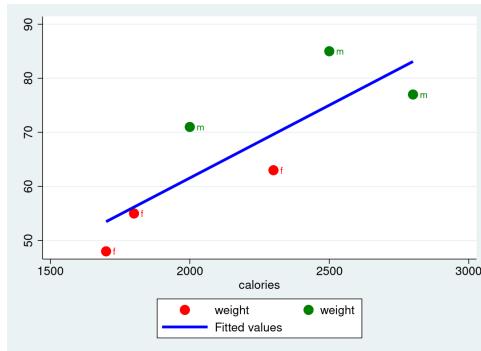
i Use `spec()` to retrieve the full column specification for this data.

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

- g) Show for all numerical variables the summary statistics including the mean, median, minimum, and the maximum.
- h) Show for all numerical variables the summary statistics including the mean and the standard deviation, **separated by male and female**. Use therefore the pipe operator.
- i) Suppose you want to analyze the general impact of average calories consumption per day on the weight. Discuss if the sample design is appropriate to draw conclusions on the population. What may cause some bias in the data? Discuss possibilities to improve the sampling and the survey, respectively.
- j) The following plot visualizes the two variables weight and calories. Discuss what can be improved in the graphical visualization.

## 9. Collection of exercises

Figure 9.2.: Weight vs. Calories



- k) Make a scatterplot matrix containing all numerical variables.
- l) Calculate the Pearson Correlation Coefficient of the two variables
  1. `calories` and `sport`
  2. `weight` and `calories`
- m) Make a scatterplot with `weight` in the y-axis and `calories` on the x-axis. Additionally, the plot should contain a linear fit and the points should be labeled with the `sex` just like in the figure shown above.
- n) Estimate the following regression specification using the OLS method: [weight\_i =  $\beta_0 + \beta_1$  calories\_i +  $\epsilon_i$ ]

Show a summary of the estimates that look like the following:

Call:

```
lm(formula = weight ~ calories, data = df)
```

Residuals:

1	2	3	4	5	6
-5.490	-1.182	-6.640	9.435	-6.099	9.976

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	7.730275	20.197867	0.383	0.7214
calories	0.026917	0.009107	2.956	0.0417 *

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8.68 on 4 degrees of freedom

Multiple R-squared: 0.6859, Adjusted R-squared: 0.6074

F-statistic: 8.735 on 1 and 4 DF, p-value: 0.04174

- o) Interpret the results. In particular, interpret how many kg the estimated weight increases—on average and ceteris paribus—if calories increase by 100 calories. Additionally, discuss the statistical properties of the estimated coefficient  $\hat{\beta}_1$  and the meaning of the **Adjusted R-squared**.

## 9. Collection of exercises

- p) OLS estimates can suffer from omitted variable bias. State the two conditions that need to be fulfilled for omitted bias to occur.
- q) Discuss potential confounding variables that may cause omitted variable bias. Given the dataset above how can some of the confounding variables be *controlled for*?

### Solution

The script uses the following functions: `aes`, `c`, `cor`, `data.frame`, `geom_point`, `geom_text`, `ggplot`, `group_by`, `lm`, `mean`, `plot`, `read_csv`, `sd`, `stat_smooth`, `summarise`, `summary`.

### R script

## 9. Collection of exercises

```
# 1
# Stephan Huber, 000, 2020-May-30

# 2
# setwd("/home/sthu/Dropbox/hsf/22-ss/dsb_bac/work/")

# 3
rm(list = ls())

# 4
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, haven)

# 5
# cross-section

# 6
sex <- c("f", "f", "f", "m", "m", "m")
age <- c(21, 19, 23, 18, 20, 61)
weight <- c(48, 55, 63, 71, 77, 85)
calories <- c(1700, 1800, 2300, 2000, 2800, 2500)
sport <- c(60, 120, 180, 60, 240, 30)
df <- data.frame(sex, age, weight, calories, sport)

# write_csv(df, file = "/home/sthu/Dropbox/hsf/exams/21-04/stuff/df.csv")
# write_csv(df, file = "/home/sthu/Dropbox/hsf/github/courses/dta/df-calories.csv")
df <- read_csv("https://raw.githubusercontent.com/hubchev/courses/main/dta/df-calories.csv")

# 7
summary(df)

# 8
df |>
  group_by(sex) |>
  summarise(
    mcal = mean(calories),
    sdcal = sd(calories),
    mweight = mean(weight),
    sdweight = sd(weight)
  )

# 9
# discussed in class

# 10
# Many things can be mentioned here such as the use of colors
# (red/blue is not a good choice for color blind people),
# the legend makes no sense as red and green both refer to \textit{sport},
# the label of 'f' and 'm' is not explained in the legend,
# rotating the labels of the y-axis would increase readability, and
# both axes do not start at zero which is hard to see.
# Also, it is a common to draw the variable you want to explain
# (here: calories) on the y-axis.

# 11
plot(df)
```

## Output of the R script

```
# 1
# Stephan Huber, 000, 2020-May-30

# 2
# setwd("/home/sthu/Dropbox/hsf/22-ss/dsb_bac/work/")

# 3
rm(list = ls())

# 4
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, haven)

# 5
# cross-section

# 6
sex <- c("f", "f", "f", "m", "m", "m")
age <- c(21, 19, 23, 18, 20, 61)
weight <- c(48, 55, 63, 71, 77, 85)
calories <- c(1700, 1800, 2300, 2000, 2800, 2500)
sport <- c(60, 120, 180, 60, 240, 30)
df <- data.frame(sex, age, weight, calories, sport)

# write_csv(df, file = "/home/sthu/Dropbox/hsf/exams/21-04/stuff/df.csv")
# write_csv(df, file = "/home/sthu/Dropbox/hsf/github/courses/dta/df-calories.csv")
df <- read_csv("https://raw.githubusercontent.com/hubchev/courses/main/dta/df-calories.csv")
```

Rows: 6 Columns: 5

-- Column specification -----  
 Delimiter: ","  
 chr (1): sex  
 dbl (4): age, weight, calories, sport

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.

```
# 7
summary(df)
```

	sex	age	weight	calories	sport
Length:	6	Min. :18.00	Min. :48.0	Min. :1700	Min. : 30
Class :	character	1st Qu.:19.25	1st Qu.:57.0	1st Qu.:1850	1st Qu.: 60
Mode :	character	Median :20.50	Median :67.0	Median :2150	Median : 90
		Mean :27.00	Mean :66.5	Mean :2183	Mean :115
		3rd Qu.:22.50	3rd Qu.:75.5	3rd Qu.:2450	3rd Qu.:165
		Max. :61.00	Max. :85.0	Max. :2800	Max. :240

## 9. Collection of exercises

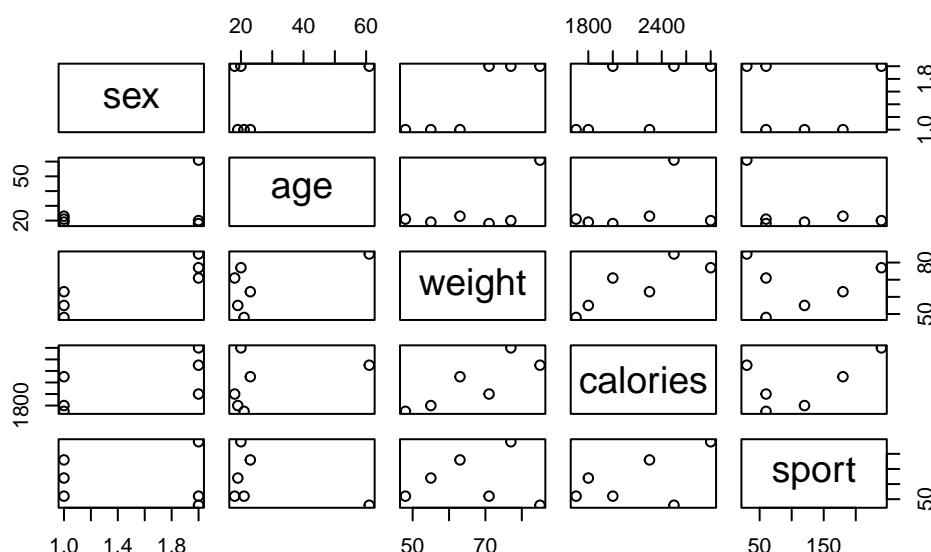
```
# 8
df |>
  group_by(sex) |>
  summarise(
    mcal = mean(calories),
    sdcal = sd(calories),
    mweight = mean(weight),
    sdweight = sd(weight)
  )
```

```
# A tibble: 2 x 5
  sex     mcal   sdcal   mweight   sdweight
  <chr>   <dbl>   <dbl>     <dbl>      <dbl>
1 f        1933.   321.     55.3       7.51
2 m        2433.   404.     77.7       7.02
```

```
# 9
# discussed in class
```

```
# 10
# Many things can be mentioned here such as the use of colors
# (red/blue is not a good choice for color blind people),
# the legend makes no sense as red and green both refer to \textit{sport},
# the label of `f` and `m` is not explained in the legend,
# rotating the labels of the y-axis would increase readability, and
# both axes do not start at zero which is hard to see.
# Also, it is a common to draw the variable you want to explain
# (here: calories) on the y-axis.
```

```
# 11
plot(df)
```



```
# 12
cor(df$calories, df$sport, method = c("pearson"))
```

9. Collection of exercises

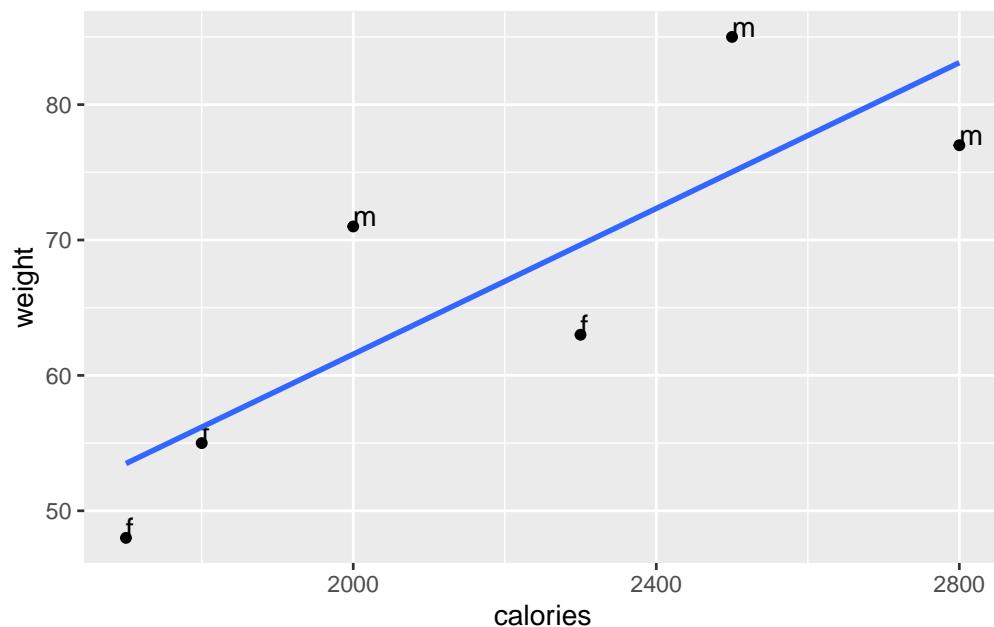
```
[1] 0.5330615
```

```
cor(df$weight, df$calories, method = c("pearson"))
```

```
[1] 0.8281972
```

```
# 13
ggplot(df, aes(x = calories, y = weight, label = sex)) +
  geom_point() +
  geom_text(hjust = 0, vjust = 0) +
  stat_smooth(formula = y ~ x, method = "lm", se = FALSE)
```

Warning: The following aesthetics were dropped during statistical transformation:  
 i This can happen when ggplot fails to infer the correct grouping structure in  
 the data.  
 i Did you forget to specify a `group` aesthetic or to convert a numerical  
 variable into a factor?



```
# 14
reg_base <- lm(weight ~ calories, data = df)
summary(reg_base)
```

Call:  
`lm(formula = weight ~ calories, data = df)`

Residuals:

1	2	3	4	5	6
-5.490	-1.182	-6.640	9.435	-6.099	9.976

```

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 7.730275  20.197867   0.383  0.7214
calories     0.026917   0.009107   2.956  0.0417 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8.68 on 4 degrees of freedom
Multiple R-squared:  0.6859,    Adjusted R-squared:  0.6074
F-statistic: 8.735 on 1 and 4 DF,  p-value: 0.04174

# 15
# 1) An increase of 100 calories (taken on average on a daily basis) is associated
# - on average and ceteris paribus - with 2.69 more of kg the participants are
# pretended to weight.
# 2) The estimated coefficient $beta_1$ is statistically significantly different to zero
# on a significance level of 5%.
# 3) About 60 % of the variation of the weight is explained by the
# estimated coefficients of the empirical model.

# 16
# For omitted variable bias to occur, the omitted variable `Z` must satisfy
# two conditions:
# 1) The omitted variable is correlated with the included regressor
# 2) The omitted variable is a determinant of the dependent variable

# 17
# discussed in class

# unload packages
suppressMessages(pacman::p_unload(tidyverse, haven))

```

## 9.15. Bundesliga

Open the script that you find [here](#) and work on the following tasks:

1. Set your working directory.
2. Clear the environment.
3. Install and load the `bundesligR` and `tidyverse`.
4. Read in the data `bundesligR` as a tibble.
5. Replace “Bor. Moenchengladbach” with “Borussia Moenchengladbach.”
6. Check for the data class.
7. View the data.
8. Glimpse on the data.

## 9. Collection of exercises

9. Show the first and last observations.
10. Show summary statistics to all variables.
11. How many teams have played in the league over the years?
12. Which teams have played Bundesliga so far?
13. How many teams have played Bundesliga?
14. How often has each team played in the Bundesliga?
15. Show summary statistics of variable `Season` only.
16. Show summary statistics of all numeric variables (`Team` is a character).
17. What is the highest number of points ever received by a team? Show only the name of the club with the highest number of points ever received.
18. Create a new tibble using `liga` removing the variable `Pts_pre_95` from the data.
19. Create a new tibble using `liga` renaming W, D, and L to Win, Draw, and Loss. Additionally rename GF, GA, GD to Goals\_shot, Goals\_received, Goal\_difference.
20. Create a new tibble using `liga` without the variable `Pts_pre_95` and only observations before the year 1996.
21. Remove the three tibbles just created from the environment.
22. Rename all variables of `liga` to lowercase and store it as `dfb`.
23. Show the winner and the runner up after the year 2010. Additionally show the points received.
24. Create a variable that counts how often a team was ranked first.
25. How often has each team played in the Bundesliga?
26. Make a ranking that shows which team has played the Bundesliga most often.
27. Add a variable to `dfb` that contains the number of appearances of a team in the league.
28. Create a number that indicates how often a team has played Bundesliga in a given year.
29. Make a ranking with the number of titles of all teams that ever won the league.
30. Create a numeric identifying variable for each team.
31. When a team is in the league, what is the probability that it wins the league?
32. Make a scatterplot with points on the y-axis and position on the x-axis.
33. Make a scatterplot with points on the y-axis and position on the x-axis. Additionally, only consider seasons with 18 teams and add lines that make clear how many points you needed to be placed in between rank 2 and 15.
34. Remove all objects from the environment except `dfb` and `liga`.
35. In Figure Figure 9.3, the ranking history of 1. FC Kaiserslautern is shown. Replicate that plot.
34. In Figure Figure 9.4, I made the graph a bit nicer. Can you spot all differences and can you guess what the dashed line and the triangles mean? How could the visualization be improved further? Replicate the plot.

Figure 9.3.: Ranking history: 1. FC Kaiserslautern

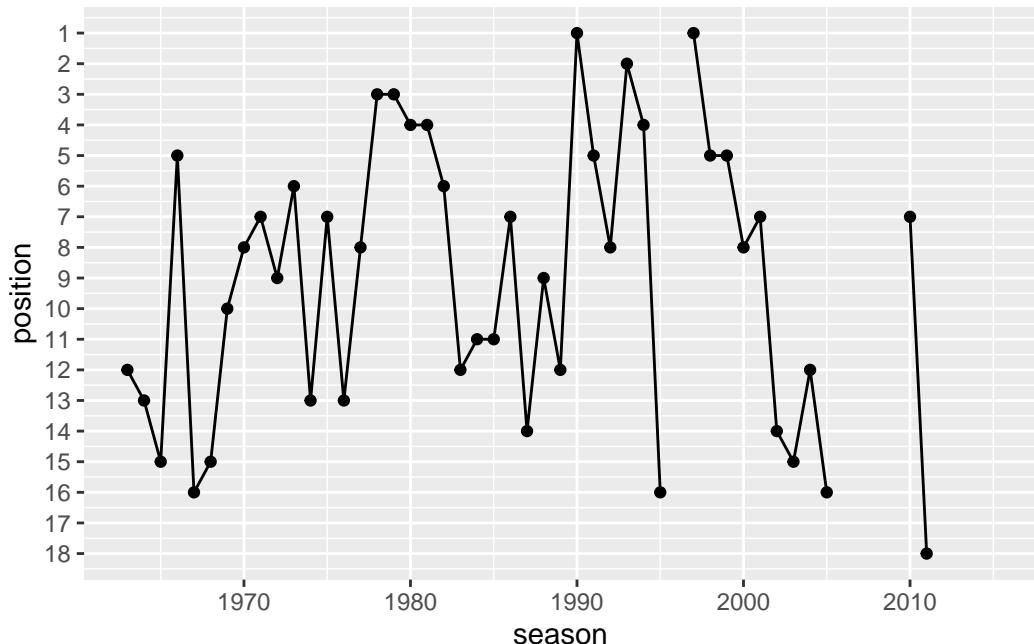
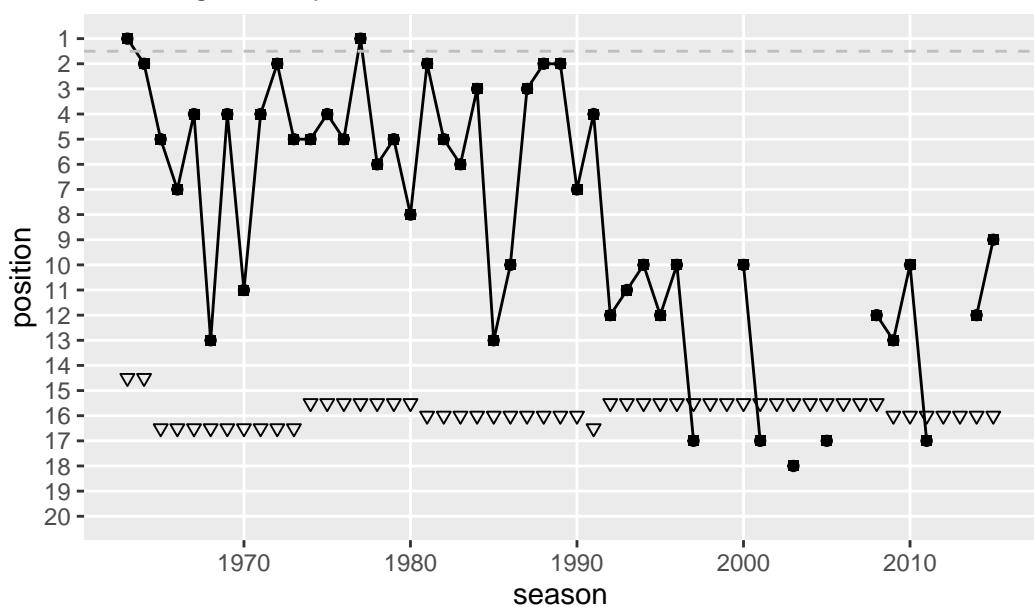


Figure 9.4.: Ranking history: 1. FC Köln  
Ranking History: 1. FC Koeln



## 9. Collection of exercises

35. Try to make the ranking history for each club ever played the league and export the graph as a png file.

### Solution

The script uses the following functions: `aes`, `arrange`, `as_tibble`, `as.numeric`, `between`, `c`, `case_when`, `class`, `complete`, `desc`, `dir.create`, `dir.exists`, `element_blank`, `facet_wrap`, `factor`, `filter`, `geom_hline`, `geom_line`, `geom_point`, `geom_vline`, `ggplot`, `ggsave`, `glimpse`, `group_by`, `head`, `ifelse`, `is.na`, `labs`, `list`, `max`, `mutate`, `n_distinct`, `paste`, `print`, `rename`, `rename_all`, `row_number`, `scale_x_continuous`, `scale_y_continuous`, `scale_y_reverse`, `select`, `seq`, `setdiff`, `slice_head`, `subset`, `sum`, `summarise`, `summary`, `table`, `tail`, `theme`, `theme_classic`, `theme_minimal`, `unique`, `unlink`, `view`, `xlim`.

### R script

## 9. Collection of exercises

```
# In dfb.R I analyze German soccer results

# set working directory
# setwd("~/Dropbox/hsf/23-ws/dsda/scripts")

# clear environment
rm(list = ls())

# (Install and) load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(
  bundesligR,
  tidyverse
)

# Read in the data as tibble
liga <- as_tibble(bundesligR)

# -----
# !!! ERRORS / ISSUES:
# "Borussia Moenchengladbach" is also entitled "Bor. Moenchengladbach"!
# Leverkusen is falsely entitled "SV Bayer 04 Leverkusen"
# Uerdingen has changed its name several times
# Stuttgarter Kickers are named differently

# How often is "Bor. Moenchengladbach" in the data?
sum(liga$Team == "Bor. Moenchengladbach")

# show the entries
liga |>
  filter(Team == "Bor. Moenchengladbach")

# Replace "Bor. Moenchengladbach" with "Borussia Moenchengladbach"
liga <- liga |>
  mutate(Team = ifelse(Team == "Bor. Moenchengladbach",
    "Borussia Moenchengladbach",
    Team
  )) |>
  mutate(Team = ifelse(Team == "SV Bayer 04 Leverkusen",
    "TSV Bayer 04 Leverkusen",
    Team
  )) |>
  mutate(Team = ifelse(Team == "FC Bayer 05 Uerdingen" |
    Team == "Bayer 05 Uerdingen",
    "KFC Uerdingen 05",
    Team
  )) |>
  mutate(Team = ifelse(Team == "SV Stuttgarter Kickers",
    "Stuttgarter Kickers",
    Team
  ))

# -----
```

```
# Check for the data class
class(liga)
```

## Output of the R script

```
# In dfb.R I analyze German soccer results

# set working directory
# setwd("~/Dropbox/hsf/23-ws/dsda/scripts")

# clear environment
rm(list = ls())

# (Install and) load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(
  bundesligR,
  tidyverse
)

# Read in the data as tibble
liga <- as_tibble(bundesligR)

# -----
# !!! ERRORS / ISSUES:
# "Borussia Moenchengladbach" is also entitled "Bor. Moenchengladbach"!
# Leverkusen is falsely entitled "SV Bayer 04 Leverkusen"
# Uerdingen has changed its name several times
# Stuttgarter Kickers are named differently

# How often is "Bor. Moenchengladbach" in the data?
sum(liga$Team == "Bor. Moenchengladbach")
```

[1] 2

```
# show the entries
liga |>
  filter(Team == "Bor. Moenchengladbach")
```

```
# A tibble: 2 x 12
  Season Position Team      Played     W     D     L     GF     GA     GD Points
  <dbl>     <dbl> <chr>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 1989       15 Bor. Moench~     34    11     8    15    37    45    -8     41
2 1976        1 Bor. Moench~     34    17    10     7    58    34     24     61
# i 1 more variable: Pts_pre_95 <dbl>
```

## 9. Collection of exercises

```
# Replace "Bor. Moenchengladbach" with "Borussia Moenchengladbach"
liga <- liga |>
  mutate(Team = ifelse(Team == "Bor. Moenchengladbach",
    "Borussia Moenchengladbach",
    Team
  )) |>
  mutate(Team = ifelse(Team == "SV Bayer 04 Leverkusen",
    "TSV Bayer 04 Leverkusen",
    Team
  )) |>
  mutate(Team = ifelse(Team == "FC Bayer 05 Uerdingen" |
    Team == "Bayer 05 Uerdingen",
    "KFC Uerdingen 05",
    Team
  )) |>
  mutate(Team = ifelse(Team == "SV Stuttgarter Kickers",
    "Stuttgarter Kickers",
    Team
  ))

# -----
```

```
# Check for the data class
class(liga)
```

```
[1] "tbl_df"      "tbl"        "data.frame"
```

```
# view data
view(liga)
```

```
# Glimpse on the data
glimpse(liga)
```

Rows: 952

Columns: 12

```
$ Season      <dbl> 2015, 2015, 2015, 2015, 2015, 2015, 2015, 2015, 2015, ~
$ Position     <dbl> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, ~
$ Team         <chr> "FC Bayern Muenchen", "Borussia Dortmund", "Bayer 04 Leverk~
$ Played       <dbl> 34, 34, 34, 34, 34, 34, 34, 34, 34, 34, 34, 34, 34, 34, 34, ~
$ W            <dbl> 28, 24, 18, 17, 15, 14, 14, 12, 10, 11, 10, 9, 10, 9, 9, 9, ~
$ D            <dbl> 4, 6, 6, 4, 7, 8, 8, 9, 13, 8, 10, 11, 8, 11, 10, 9, 6, 4, ~
$ L            <dbl> 2, 4, 10, 13, 12, 12, 12, 13, 11, 15, 14, 14, 16, 14, 15, 1~
$ GF           <dbl> 80, 82, 56, 67, 51, 46, 42, 47, 38, 40, 33, 42, 50, 38, 39, ~
$ GA           <dbl> 17, 34, 40, 50, 49, 42, 42, 49, 42, 46, 42, 52, 65, 53, 54, ~
$ GD           <dbl> 63, 48, 16, 17, 2, 4, 0, -2, -4, -6, -9, -10, -15, -15, -15~
$ Points        <dbl> 88, 78, 60, 55, 52, 50, 50, 45, 43, 41, 40, 38, 38, 38, 37, ~
$ Pts_pre_95   <dbl> 60, 54, 42, 38, 37, 36, 33, 33, 30, 30, 29, 28, 29, 28, ~
```

```
# first and last observations
head(liga)
```

## 9. Collection of exercises

```
# A tibble: 6 x 12
  Season Position Team      Played     W      D      L      GF      GA      GD Points
  <dbl>    <dbl> <chr>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 2015        1 FC Bayern M~     34     28     4     2     80     17     63     88
2 2015        2 Borussia Do~     34     24     6     4     82     34     48     78
3 2015        3 Bayer 04 Le~     34     18     6    10     56     40     16     60
4 2015        4 Borussia Mo~     34     17     4    13     67     50     17     55
5 2015        5 FC Schalke ~     34     15     7    12     51     49      2     52
6 2015        6 1. FSV Main~     34     14     8    12     46     42      4     50
# i 1 more variable: Pts_pre_95 <dbl>
```

```
tail(liga)
```

```
# A tibble: 6 x 12
  Season Position Team      Played     W      D      L      GF      GA      GD Points
  <dbl>    <dbl> <chr>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 1963        11 Eintracht B~     30     11     6    13     36     49    -13     39
2 1963        12 1. FC Kaise~     30     10     6    14     48     69    -21     36
3 1963        13 Karlsruher ~     30      8     8    14     42     55    -13     32
4 1963        14 Hertha BSC     30      9     6    15     45     65    -20     33
5 1963        15 Preussen Mu~     30      7     9    14     34     52    -18     30
6 1963        16 1. FC Saarb~     30      6     5    19     44     72    -28     23
# i 1 more variable: Pts_pre_95 <dbl>
```

```
# summary statistics
summary(liga)
```

Season	Position	Team	Played
Min. :1963	Min. : 1.000	Length:952	Min. :30.00
1st Qu.:1976	1st Qu.: 5.000	Class :character	1st Qu.:34.00
Median :1989	Median : 9.000	Mode :character	Median :34.00
Mean :1989	Mean : 9.486		Mean :33.95
3rd Qu.:2002	3rd Qu.:14.000		3rd Qu.:34.00
Max. :2015	Max. :20.000		Max. :38.00
W	D	L	GF
Min. : 2.00	Min. : 2.000	Min. : 1.00	Min. : 15.00
1st Qu.: 9.75	1st Qu.: 7.000	1st Qu.:10.00	1st Qu.: 42.00
Median :12.00	Median : 9.000	Median :13.00	Median : 50.00
Mean :12.61	Mean : 8.733	Mean :12.61	Mean : 52.01
3rd Qu.:15.00	3rd Qu.:11.000	3rd Qu.:15.00	3rd Qu.: 61.00
Max. :29.00	Max. :18.000	Max. :28.00	Max. :101.00
GA	GD	Points	Pts_pre_95
Min. :10.0	Min. : -60.0000	Min. :10.00	Min. : 8.00
1st Qu.:43.0	1st Qu.: -13.0000	1st Qu.:38.00	1st Qu.:29.00
Median :51.0	Median : -2.0000	Median :44.00	Median :33.00
Mean :51.7	Mean : 0.3015	Mean :46.56	Mean :33.95
3rd Qu.:60.0	3rd Qu.: 13.0000	3rd Qu.:55.00	3rd Qu.:39.00
Max. :93.0	Max. : 80.0000	Max. :91.00	Max. :62.00

## 9. Collection of exercises

```
# How many teams have played in the league over the years?


```

```
1963 1964 1965 1966 1967 1968 1969 1970 1971 1972 1973 1974 1975 1976 1977 1978
   16    16    18    18    18    18    18    18    18    18    18    18    18    18    18    18
1979 1980 1981 1982 1983 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994
   18    18    18    18    18    18    18    18    18    18    18    18    18    20    18    18
1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010
   18    18    18    18    18    18    18    18    18    18    18    18    18    18    18    18
2011 2012 2013 2014 2015
   18    18    18    18    18
```

```
# Which teams have played Bundesliga
unique(liga$Team)
```

```
[1] "FC Bayern Muenchen"      "Borussia Dortmund"
[3] "Bayer 04 Leverkusen"     "Borussia Moenchengladbach"
[5] "FC Schalke 04"          "1. FSV Mainz 05"
[7] "Hertha BSC"             "VfL Wolfsburg"
[9] "1. FC Koeln"            "Hamburger SV"
[11] "FC Ingolstadt 04"       "FC Augsburg"
[13] "Werder Bremen"          "SV Darmstadt 98"
[15] "TSG 1899 Hoffenheim"    "Eintracht Frankfurt"
[17] "VfB Stuttgart"          "Hannover 96"
[19] "SC Freiburg"            "SC Paderborn 07"
[21] "1. FC Nuernberg"        "Eintracht Braunschweig"
[23] "Fortuna Duesseldorf"   "SpVgg Greuther Fuerth"
[25] "1. FC Kaiserslautern"   "FC St. Pauli"
[27] "VfL Bochum"             "Energie Cottbus"
[29] "Karlsruher SC"          "Arminia Bielefeld"
[31] "Hansa Rostock"          "MSV Duisburg"
[33] "Alemannia Aachen"        "TSV 1860 Muenchen"
[35] "SpVgg Unterhaching"     "SSV Ulm 1846"
[37] "KFC Uerdingen 05"        "Dynamo Dresden"
[39] "SG Wattenscheid 09"      "VfB Leipzig"
[41] "1. FC Saarbruecken"     "TSV Bayer 04 Leverkusen"
[43] "SV Werder Bremen"        "1. FC Dynamo Dresden"
[45] "Stuttgarter Kickers"     "FC Hansa Rostock"
[47] "SV Waldhof Mannheim"    "FC 08 Homburg"
[49] "FC Homburg"              "Blau-Weiss 90 Berlin"
[51] "Kickers Offenbach"       "Tennis Borussia Berlin"
[53] "Rot-Weiss Essen"          "Wuppertaler SV"
[55] "SC Fortuna Koeln"        "Rot-Weiss Oberhausen"
[57] "SC Rot-Weiss Oberhausen" "Borussia Neunkirchen"
[59] "Meidericher SV"          "SC Tasmania 1900 Berlin"
[61] "Preussen Muenster"
```

```
# How many teams have played Bundesliga
n_distinct(liga$Team)
```

9. Collection of exercises

[1] 61

```
# How often has each team played in the Bundesliga
table(liga$Team)
```

1. FC Dynamo Dresden	1	1. FC Kaiserslautern	44	1. FC Koeln	45
1. FC Nuernberg	32	1. FC Saarbruecken	5	1. FSV Mainz 05	10
Alemannia Aachen	4	Arminia Bielefeld	17	Bayer 04 Leverkusen	30
Blau-Weiss 90 Berlin	1	Borussia Dortmund	49	Borussia Moenchengladbach	48
Borussia Neunkirchen	3	Dynamo Dresden	3	Eintracht Braunschweig	21
Eintracht Frankfurt	47	Energie Cottbus	6	FC 08 Homburg	2
FC Augsburg	5	FC Bayern Muenchen	51	FC Hansa Rostock	1
FC Homburg	1	FC Ingolstadt 04	1	FC Schalke 04	48
FC St. Pauli	8	Fortuna Duesseldorf	23	Hamburger SV	53
Hannover 96	28	Hansa Rostock	11	Hertha BSC	33
Karlsruher SC	24	KFC Uerdingen 05	14	Kickers Offenbach	7
Meidericher SV	3	MSV Duisburg	25	Preussen Muenster	1
Rot-Weiss Essen	7	Rot-Weiss Oberhausen	3	SC Fortuna Koeln	1
SC Freiburg	16	SC Paderborn 07	1	SC Rot-Weiss Oberhausen	1
SC Tasmania 1900 Berlin	1	SG Wattenscheid 09	4	SpVgg Greuther Fuerth	1
SpVgg Unterhaching	2	SSV Ulm 1846	1	Stuttgarter Kickers	2
SV Darmstadt 98	3	SV Waldhof Mannheim	7	SV Werder Bremen	1
Tennis Borussia Berlin	2	TSG 1899 Hoffenheim	8	TSV 1860 Muenchen	20
TSV Bayer 04 Leverkusen	7	VfB Leipzig	1	VfB Stuttgart	51
VfL Bochum	34	VfL Wolfsburg	19	Werder Bremen	51
Wuppertaler SV	3				

## 9. Collection of exercises

```
# summary of variable Season only
summary(liga$Season)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1963	1976	1989	1989	2002	2015

```
# summary of numeric of variables (Team is a character)
liga |>
  select(Season, Position, Played, W, D, L, GF, GA, GD, Points, Pts_pre_95) |>
  summary()
```

Season	Position	Played	W	
Min. :1963	Min. : 1.000	Min. :30.00	Min. : 2.00	
1st Qu.:1976	1st Qu.: 5.000	1st Qu.:34.00	1st Qu.: 9.75	
Median :1989	Median : 9.000	Median :34.00	Median :12.00	
Mean :1989	Mean : 9.486	Mean :33.95	Mean :12.61	
3rd Qu.:2002	3rd Qu.:14.000	3rd Qu.:34.00	3rd Qu.:15.00	
Max. :2015	Max. :20.000	Max. :38.00	Max. :29.00	
	D	L	GF	
	Min. : 2.000	Min. : 1.00	Min. : 15.00	Min. :10.0
	1st Qu.: 7.000	1st Qu.:10.00	1st Qu.: 42.00	1st Qu.:43.0
	Median : 9.000	Median :13.00	Median : 50.00	Median :51.0
	Mean : 8.733	Mean :12.61	Mean : 52.01	Mean :51.7
	3rd Qu.:11.000	3rd Qu.:15.00	3rd Qu.: 61.00	3rd Qu.:60.0
	Max. :18.000	Max. :28.00	Max. :101.00	Max. :93.0
	GD	Points	Pts_pre_95	
	Min. : -60.0000	Min. :10.00	Min. : 8.00	
	1st Qu.: -13.0000	1st Qu.:38.00	1st Qu.:29.00	
	Median : -2.0000	Median :44.00	Median :33.00	
	Mean : 0.3015	Mean :46.56	Mean :33.95	
	3rd Qu.: 13.0000	3rd Qu.:55.00	3rd Qu.:39.00	
	Max. : 80.0000	Max. :91.00	Max. :62.00	

```
# shorter alternative
liga |>
  select(Season, Position, Played:Pts_pre_95) |>
  summary()
```

Season	Position	Played	W	
Min. :1963	Min. : 1.000	Min. :30.00	Min. : 2.00	
1st Qu.:1976	1st Qu.: 5.000	1st Qu.:34.00	1st Qu.: 9.75	
Median :1989	Median : 9.000	Median :34.00	Median :12.00	
Mean :1989	Mean : 9.486	Mean :33.95	Mean :12.61	
3rd Qu.:2002	3rd Qu.:14.000	3rd Qu.:34.00	3rd Qu.:15.00	
Max. :2015	Max. :20.000	Max. :38.00	Max. :29.00	
	D	L	GF	
	Min. : 2.000	Min. : 1.00	Min. : 15.00	Min. :10.0
	1st Qu.: 7.000	1st Qu.:10.00	1st Qu.: 42.00	1st Qu.:43.0

## 9. Collection of exercises

```

Median : 9.000   Median :13.00   Median : 50.00   Median :51.0
Mean   : 8.733   Mean   :12.61   Mean   : 52.01   Mean   :51.7
3rd Qu.:11.000   3rd Qu.:15.00   3rd Qu.: 61.00   3rd Qu.:60.0
Max.   :18.000   Max.   :28.00   Max.   :101.00  Max.   :93.0
      GD          Points        Pts_pre_95
Min.  :-60.0000  Min.   :10.00   Min.   : 8.00
1st Qu.:-13.0000 1st Qu.:38.00  1st Qu.:29.00
Median : -2.0000  Median :44.00   Median :33.00
Mean   : 0.3015   Mean   :46.56   Mean   :33.95
3rd Qu.: 13.0000 3rd Qu.:55.00  3rd Qu.:39.00
Max.   : 80.0000  Max.   :91.00   Max.   :62.00

```

```
# shortest alternative
liga |>
  select(-Team) |>
  filter(Season == 1999 | Season == 2010) |>
  summary()
```

Season	Position	Played	W	D
Min. :1999	Min. : 1.0	Min. :34	Min. : 4.00	Min. : 3.000
1st Qu.:1999	1st Qu.: 5.0	1st Qu.:34	1st Qu.: 9.75	1st Qu.: 6.000
Median :2004	Median : 9.5	Median :34	Median :12.00	Median : 8.000
Mean   :2004	Mean   : 9.5	Mean   :34	Mean   :12.83	Mean   : 8.333
3rd Qu.:2010	3rd Qu.:14.0	3rd Qu.:34	3rd Qu.:14.25	3rd Qu.:10.250
Max.   :2010	Max.   :18.0	Max.   :34	Max.   :23.00	Max.   :15.000
L	GF	GA	GD	
Min. : 3.00	Min. :31.00	Min. :22.00	Min. : -34.00	
1st Qu.:10.75	1st Qu.:41.00	1st Qu.:44.00	1st Qu.: -10.25	
Median :13.00	Median :47.00	Median :48.50	Median : -3.00	
Mean   :12.83	Mean   :49.42	Mean   :49.42	Mean   : 0.00	
3rd Qu.:16.00	3rd Qu.:54.25	3rd Qu.:59.00	3rd Qu.: 4.75	
Max.   :21.00	Max.   :81.00	Max.   :71.00	Max.   : 45.00	
Points	Pts_pre_95			
Min. :22.00	Min. :18.00			
1st Qu.:39.75	1st Qu.:29.75			
Median :44.00	Median :32.00			
Mean   :46.83	Mean   :34.00			
3rd Qu.:50.75	3rd Qu.:37.50			
Max.   :75.00	Max.   :52.00			

```
# Most points ever received by a team
liga |>
  filter(Points == max(Points))
```

```
# A tibble: 1 x 12
  Season Position Team       Played     W     D     L     GF     GA     GD Points
  <dbl>    <dbl> <chr>     <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 2012        1 FC Bayern M~     34    29     4     1    98    18     80     91
# i 1 more variable: Pts_pre_95 <dbl>
```

## 9. Collection of exercises

```

# Show only the team name
liga |>
  filter(Points == max(Points)) |>
  select(Team) |>
  print()

# A tibble: 1 x 1
  Team
  <chr>
1 FC Bayern Muenchen

# remove the variable `Pts_pre_95` from the data
liga_post95 <- liga |>
  select(-Pts_pre_95)

# rename W, D, and L to Win, Draw, and Loss
# additionally rename GF, GA, GD to Goals_shot, Goals_received, Goal_difference
liga_longnames <- liga |>
  rename(Win = W, Draw = D, Loss = L) |>
  rename(Goals_shot = GF, Goals_received = GA, Goal_difference = GD)

# Remove the variable `Pts_pre_95` from `liga`
# additionally remove all observations before the year 1996
liga_no3point <- liga |>
  select(-Pts_pre_95) |>
  filter(Season >= 1996)

# Remove the objects liga_post95, liga_longnames, and liga_no3point from the environment
rm(liga_post95, liga_longnames, liga_no3point)

# Rename all variables of `liga` to lower cases and store it as `dfb`
dfb <- liga |>
  rename_all(tolower)

# Show the winner and the runner up after 2010
# additionally show the points
dfb |>
  filter(season > 2010) |>
  group_by(season) |>
  arrange(desc(points)) |>
  slice_head(n = 2) |>
  select(team, points, position)

# A tibble: 10 x 4
# Groups:   season [5]
  season team              points position
  <dbl> <chr>            <dbl>     <dbl>
1 2011 Borussia Dortmund    81        1
2 2011 FC Bayern Muenchen  73        2

```

## 9. Collection of exercises

3	2012 FC Bayern Muenchen	91	1
4	2012 Borussia Dortmund	66	2
5	2013 FC Bayern Muenchen	90	1
6	2013 Borussia Dortmund	71	2
7	2014 FC Bayern Muenchen	79	1
8	2014 VfL Wolfsburg	69	2
9	2015 FC Bayern Muenchen	88	1
10	2015 Borussia Dortmund	78	2

```
# Create a variable that counts how often a team was ranked first
dfb <- dfb |>
  group_by(team) |>
  mutate(meister_count = sum(position == 1))

# How often has each team played in the Bundesliga
table(liga$Team)
```

1.	FC Dynamo Dresden	1	1.	FC Kaiserslautern	44	1.	FC Koeln	45
	1. FC Nuernberg	32		1. FC Saarbruecken	5		1. FSV Mainz 05	10
	Alemannia Aachen	4		Arminia Bielefeld	17		Bayer 04 Leverkusen	30
	Blau-Weiss 90 Berlin	1		Borussia Dortmund	49		Borussia Moenchengladbach	48
	Borussia Neunkirchen	3		Dynamo Dresden	3		Eintracht Braunschweig	21
	Eintracht Frankfurt	47		Energie Cottbus	6		FC 08 Homburg	2
	FC Augsburg	5		FC Bayern Muenchen	51		FC Hansa Rostock	1
	FC Homburg	1		FC Ingolstadt 04	1		FC Schalke 04	48
	FC St. Pauli	8		Fortuna Duesseldorf	23		Hamburger SV	53
	Hannover 96	28		Hansa Rostock	11		Hertha BSC	33
	Karlsruher SC	24		KFC Uerdingen 05	14		Kickers Offenbach	7
	Meidericher SV	3		MSV Duisburg	25		Preussen Muenster	1
	Rot-Weiss Essen	7		Rot-Weiss Oberhausen	3		SC Fortuna Koeln	1
	SC Freiburg	16		SC Paderborn 07	1		SC Rot-Weiss Oberhausen	1
	SC Tasmania 1900 Berlin	1		SG Wattenscheid 09	4		SpVgg Greuther Fuerth	1
	SpVgg Unterhaching	2		SSV Ulm 1846	1		Stuttgarter Kickers	2

## 9. Collection of exercises

SV Darmstadt	98	SV Waldhof Mannheim	7	SV Werder Bremen	
	3		7		1
Tennis Borussia Berlin	2	TSG 1899 Hoffenheim	8	TSV 1860 Muenchen	20
TSV Bayer 04 Leverkusen	7	VfB Leipzig	1	VfB Stuttgart	51
VfL Bochum	34	VfL Wolfsburg	19	Werder Bremen	
					51
Wuppertaler SV	3				

```
# Make a ranking
dfb |>
  group_by(team) |>
  summarise(appearances = n_distinct(season)) |>
  arrange(desc(appearances)) |>
  print(n = Inf)
```

```
# A tibble: 61 x 2
  team                  appearances
  <chr>                    <int>
1 Hamburger SV            53
2 FC Bayern Muenchen     51
3 VfB Stuttgart            51
4 Werder Bremen            51
5 Borussia Dortmund        49
6 Borussia Moenchengladbach 48
7 FC Schalke 04            48
8 Eintracht Frankfurt      47
9 1. FC Koeln                45
10 1. FC Kaiserslautern       44
11 VfL Bochum                34
12 Hertha BSC                33
13 1. FC Nuernberg           32
14 Bayer 04 Leverkusen       30
15 Hannover 96                 28
16 MSV Duisburg                25
17 Karlsruher SC                24
18 Fortuna Duesseldorf        23
19 Eintracht Braunschweig       21
20 TSV 1860 Muenchen           20
21 VfL Wolfsburg                19
22 Arminia Bielefeld             17
23 SC Freiburg                  16
24 KFC Uerdingen 05              14
25 Hansa Rostock                  11
26 1. FSV Mainz 05                10
27 FC St. Pauli                  8
28 TSG 1899 Hoffenheim               8
29 Kickers Offenbach                7
```

## 9. Collection of exercises

30 Rot-Weiss Essen	7
31 SV Waldhof Mannheim	7
32 TSV Bayer 04 Leverkusen	7
33 Energie Cottbus	6
34 1. FC Saarbruecken	5
35 FC Augsburg	5
36 Alemannia Aachen	4
37 SG Wattenscheid 09	4
38 Borussia Neunkirchen	3
39 Dynamo Dresden	3
40 Meidericher SV	3
41 Rot-Weiss Oberhausen	3
42 SV Darmstadt 98	3
43 Wuppertaler SV	3
44 FC 08 Homburg	2
45 SpVgg Unterhaching	2
46 Stuttgarter Kickers	2
47 Tennis Borussia Berlin	2
48 1. FC Dynamo Dresden	1
49 Blau-Weiss 90 Berlin	1
50 FC Hansa Rostock	1
51 FC Homburg	1
52 FC Ingolstadt 04	1
53 Preussen Muenster	1
54 SC Fortuna Koeln	1
55 SC Paderborn 07	1
56 SC Rot-Weiss Oberhausen	1
57 SC Tasmania 1900 Berlin	1
58 SSV Ulm 1846	1
59 SV Werder Bremen	1
60 SpVgg Greuther Fuerth	1
61 VfB Leipzig	1

```
# Add a variable to `dfb` that contains the number of appearances of a team in the league
dfb <- dfb |>
  group_by(team) |>
  mutate(appearances = n_distinct(season))

# create a number that indicates how often a team has played Bundesliga in a given year
dfb <- dfb |>
  arrange(team, season) |>
  group_by(team) |>
  mutate(team_in_liga_count = row_number())

# Make a ranking with the number of titles of all teams that ever won the league
dfb |>
  filter(team_in_liga_count == 1) |>
  filter(meister_count != 0) |>
  arrange(desc(meister_count)) |>
  select(meister_count, team)
```

## 9. Collection of exercises

```
# A tibble: 12 x 2
# Groups:   team [12]
  meister_count team
     <int> <chr>
1          25 FC Bayern Muenchen
2           5 Borussia Dortmund
3           5 Borussia Moenchengladbach
4           4 Werder Bremen
5           3 Hamburger SV
6           3 VfB Stuttgart
7          2 1. FC Kaiserslautern
8          2 1. FC Koeln
9          1 1. FC Nuernberg
10         1 Eintracht Braunschweig
11         1 TSV 1860 Muenchen
12         1 VfL Wolfsburg

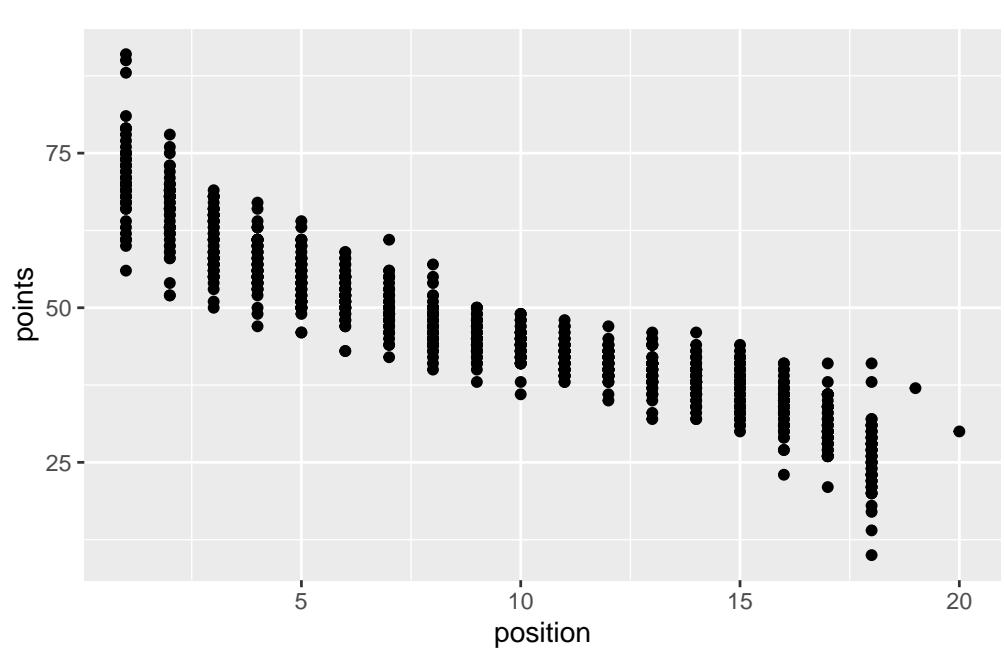
# Create a numeric identifying variable for each team
dfb_teamid <- dfb |>
  mutate(team_id = as.numeric(factor(team)))

# When a team is in the league, what is the probability that it wins the league
dfb |>
  filter(team_in_liga_count == 1) |>
  mutate(prob_win = meister_count / appearances) |>
  filter(prob_win > 0) |>
  arrange(desc(prob_win)) |>
  select(meister_count, prob_win, team)

# A tibble: 12 x 3
# Groups:   team [12]
  meister_count prob_win team
     <int>     <dbl> <chr>
1          25    0.490 FC Bayern Muenchen
2           5    0.104 Borussia Moenchengladbach
3           5    0.102 Borussia Dortmund
4           4    0.0784 Werder Bremen
5           3    0.0588 VfB Stuttgart
6           3    0.0566 Hamburger SV
7           1    0.0526 VfL Wolfsburg
8           1    0.05    TSV 1860 Muenchen
9           1    0.0476 Eintracht Braunschweig
10          2    0.0455 1. FC Kaiserslautern
11          2    0.0444 1. FC Koeln
12          1    0.0312 1. FC Nuernberg

# make a scatterplot with points on the y-axis and position on the x-axis
ggplot(dfb, aes(x = position, y = points)) +
  geom_point()
```

*9. Collection of exercises*



## 9. Collection of exercises

```
# Make a scatterplot with points on the y-axis and position on the x-axis.
# Additionally, only consider seasons with 18 teams and
# add lines that make clear how many points you needed to be placed
# in between rank 2 and 15.

dfb_18 <- dfb |>
  group_by(season) |>
  mutate(teams_in_league = n_distinct(team)) |>
  filter(teams_in_league == 18)

h_1 <- dfb_18 |>
  filter(position == 16) |>
  mutate(ma = max(points))

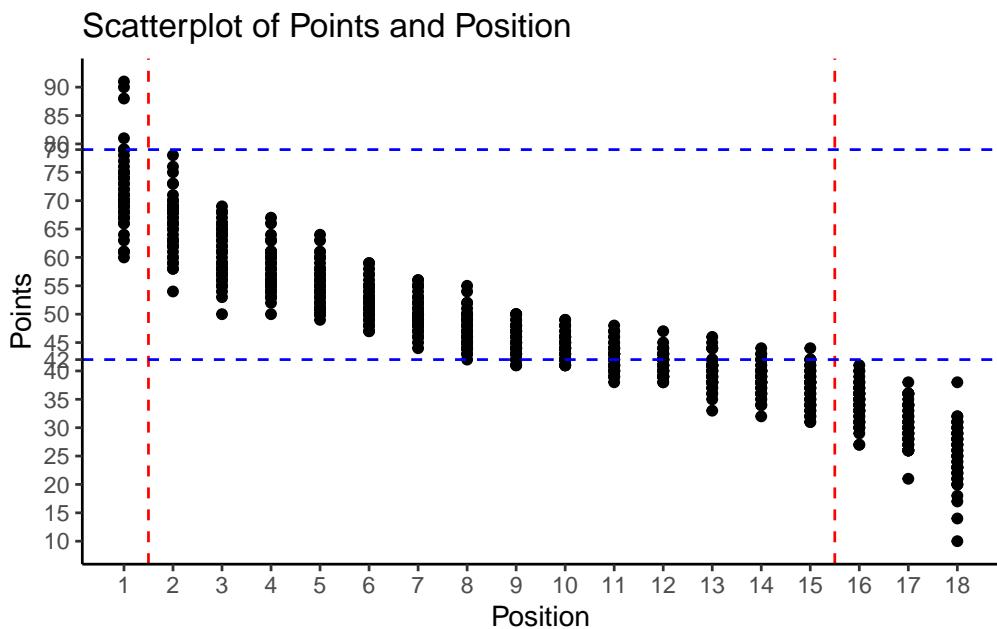
max_points_rank_16 <- max(h_1$ma) + 1

h_2 <- dfb_18 |>
  filter(position == 2) |>
  mutate(mb = max(points))

min_points_rank_2 <- max(h_2$mb) + 1

dfb_18 <- dfb_18 |>
  mutate(season_category = case_when(
    season < 1970 ~ 1,
    between(season, 1970, 1979) ~ 2,
    between(season, 1980, 1989) ~ 3,
    between(season, 1990, 1999) ~ 4,
    between(season, 2000, 2009) ~ 5,
    between(season, 2010, 2019) ~ 6,
    TRUE ~ 7 # Adjust this line based on the actual range of your data
  ))

ggplot(dfb_18, aes(x = position, y = points)) +
  geom_point() +
  labs(
    title = "Scatterplot of Points and Position",
    x = "Position",
    y = "Points"
  ) +
  geom_vline(xintercept = c(1.5, 15.5), linetype = "dashed", color = "red") +
  geom_hline(yintercept = max_points_rank_16, linetype = "dashed", color = "blue") +
  geom_hline(yintercept = min_points_rank_2, linetype = "dashed", color = "blue") +
  scale_y_continuous(breaks = c(min_points_rank_2, max_points_rank_16, seq(0, max(dfb_18$points), by = 1))) +
  scale_x_continuous(breaks = c(seq(0, max(dfb_18$points), by = 1))) +
  theme_classic()
```



```
# Remove all objects except liga and dfb
rm(list = setdiff(ls(), c("liga", "dfb")))

# Rank "1. FC Kaiserslautern" over time
dfb_bal <- dfb |>
  select(season, team, position) |>
  as_tibble() |>
  complete(season, team)

table(dfb_bal$team)
```

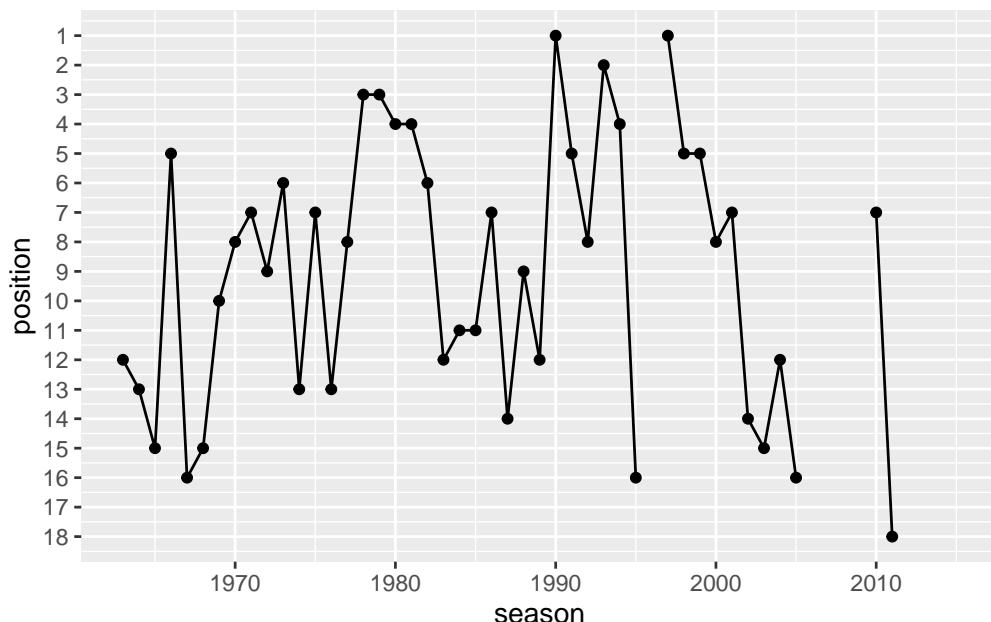
1. FC Dynamo Dresden	53	1. FC Kaiserslautern	53	1. FC Koeln	53
1. FC Nuernberg	53	1. FC Saarbruecken	53	1. FSV Mainz 05	53
Alemannia Aachen	53	Arminia Bielefeld	53	Bayer 04 Leverkusen	53
Blau-Weiss 90 Berlin	53	Borussia Dortmund	53	Borussia Moenchengladbach	53
Borussia Neunkirchen	53	Dynamo Dresden	53	Eintracht Braunschweig	53
Eintracht Frankfurt	53	Energie Cottbus	53	FC 08 Homburg	53
FC Augsburg	53	FC Bayern Muenchen	53	FC Hansa Rostock	53
FC Homburg	53	FC Ingolstadt 04	53	FC Schalke 04	53
FC St. Pauli	53	Fortuna Duesseldorf	53	Hamburger SV	53
Hannover 96	53	Hansa Rostock	53	Hertha BSC	53

## 9. Collection of exercises

Karlsruher SC		KFC Uerdingen 05		Kickers Offenbach
53		53		53
Meidericher SV		MSV Duisburg		Preussen Muenster
53		53		53
Rot-Weiss Essen		Rot-Weiss Oberhausen		SC Fortuna Koeln
53		53		53
SC Freiburg		SC Paderborn 07		SC Rot-Weiss Oberhausen
53		53		53
SC Tasmania 1900 Berlin		SG Wattenscheid 09		SpVgg Greuther Fuerth
53		53		53
SpVgg Unterhaching		SSV Ulm 1846		Stuttgarter Kickers
53		53		53
SV Darmstadt 98		SV Waldhof Mannheim		SV Werder Bremen
53		53		53
Tennis Borussia Berlin		TSG 1899 Hoffenheim		TSV 1860 Muenchen
53		53		53
TSV Bayer 04 Leverkusen		VfB Leipzig		VfB Stuttgart
53		53		53
VfL Bochum		VfL Wolfsburg		Werder Bremen
53		53		53
Wuppertaler SV				
53				

```
dfb_fck <- dfb_bal |>
  filter(team == "1. FC Kaiserslautern")

ggplot(dfb_fck, aes(x = season, y = position)) +
  geom_point() +
  geom_line() +
  scale_y_reverse(breaks = seq(1, 18, by = 1))
```



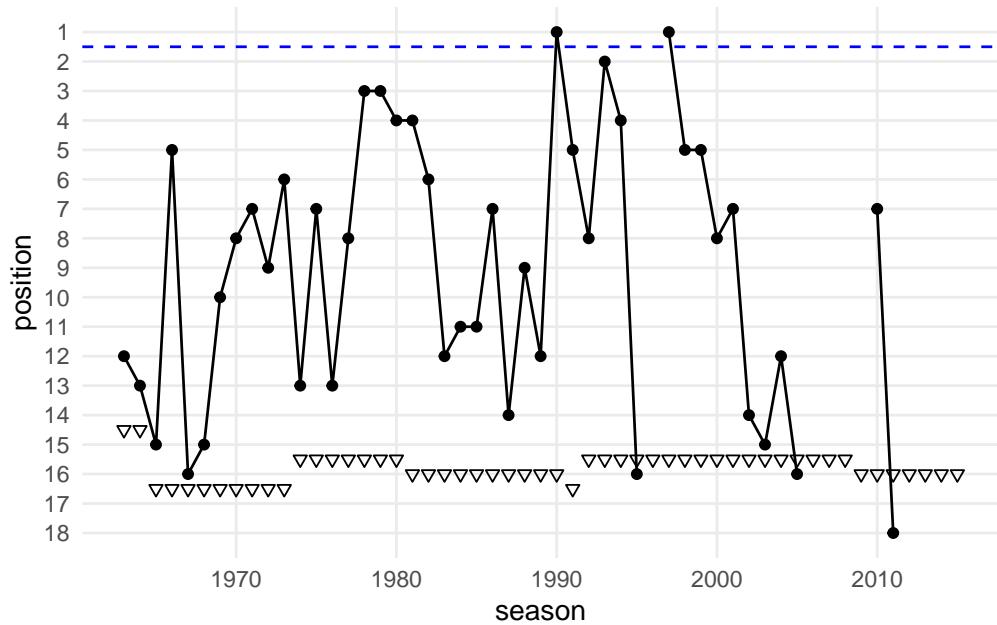
## 9. Collection of exercises

```
# Make the plot nice
```

```
# consider different rules for having to leave the league:
```

```
dfb_fck <- dfb_fck |>
  mutate(godown = ifelse(season <= 1964, 14.5, NA)) |>
  mutate(godown = ifelse(season > 1964 & season <= 1973, 16.5, godown)) |>
  mutate(godown = ifelse(season > 1973 & season <= 1980, 15.5, godown)) |>
  mutate(godown = ifelse(season > 1980 & season <= 1990, 16, godown)) |>
  mutate(godown = ifelse(season == 1991, 16.5, godown)) |>
  mutate(godown = ifelse(season > 1991 & season <= 2008, 15.5, godown)) |>
  mutate(godown = ifelse(season > 2008, 16, godown))
```

```
ggplot(dfb_fck, aes(x = season)) +
  geom_point(aes(y = position)) +
  geom_line(aes(y = position)) +
  geom_point(aes(y = godown), shape = 25) +
  scale_y_reverse(breaks = seq(1, 18, by = 1)) +
  theme_minimal() +
  theme(panel.grid.minor = element_blank()) +
  geom_hline(yintercept = 1.5, linetype = "dashed", color = "blue")
```



## 9. Collection of exercises

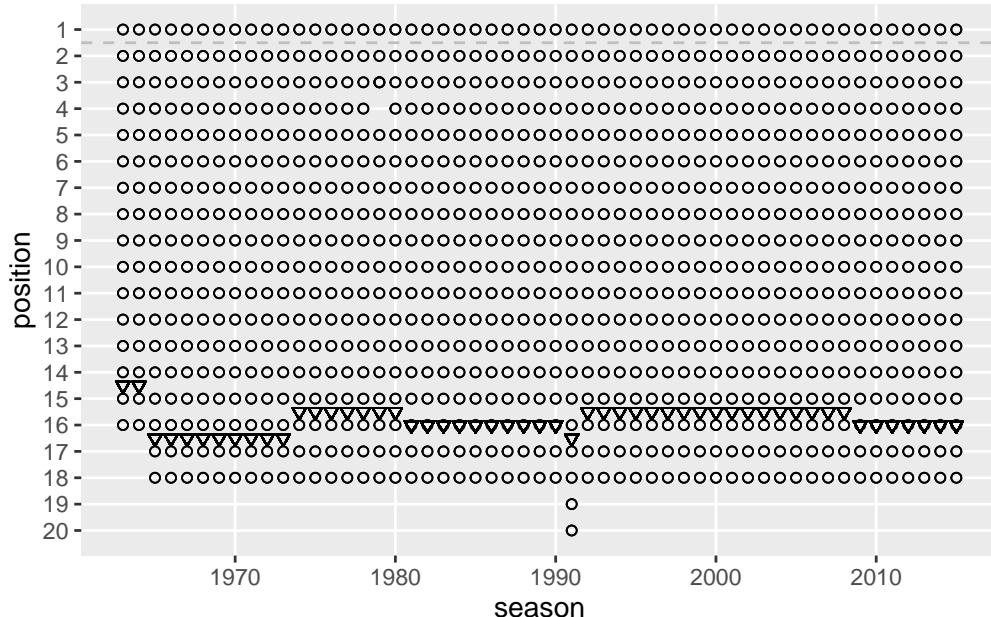
```

dfb_bal <- dfb_bal |>
  mutate(godown = ifelse(season <= 1964, 14.5, NA)) |>
  mutate(godown = ifelse(season > 1964 & season <= 1973, 16.5, godown)) |>
  mutate(godown = ifelse(season > 1973 & season <= 1980, 15.5, godown)) |>
  mutate(godown = ifelse(season > 1980 & season <= 1990, 16, godown)) |>
  mutate(godown = ifelse(season == 1991, 16.5, godown)) |>
  mutate(godown = ifelse(season > 1991 & season <= 2008, 15.5, godown)) |>
  mutate(godown = ifelse(season > 2008, 16, godown)) |>
  mutate(inliga = ifelse(is.na(position), 0, 1))

rank_plot <- ggplot(dfb_bal, aes(x = season)) +
  geom_point(aes(y = position), shape = 1) +
  # geom_line(aes(y = position)) +
  geom_point(aes(y = godown), shape = 25) +
  scale_y_reverse(breaks = seq(1, 20, by = 1), limits = c(20, 1)) +
  xlim(1963, 2015) +
  theme(panel.grid.minor = element_blank()) +
  geom_hline(yintercept = 1.5, linetype = "dashed", color = "gray") +
  geom_point(aes(y = position), shape = 1)

rank_plot

```



```

# !--> in 1979 is a gap! Error?
# No. Reason: two clubs shared the third place.

rank_plot +
  facet_wrap(~team)

```

## 9. Collection of exercises

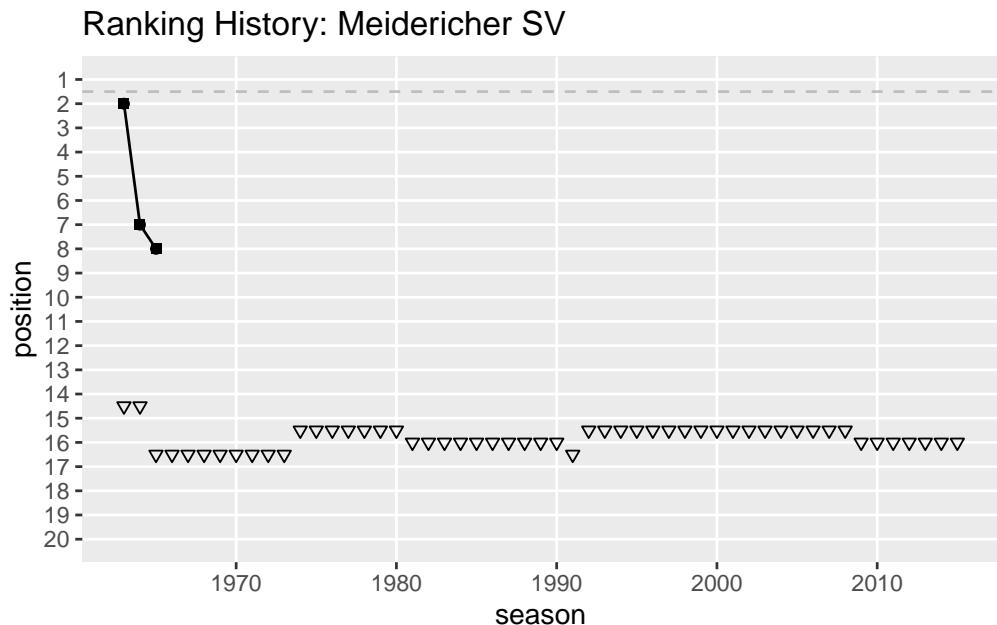


```
# Create "test" directory if it doesn't already exist
if (!dir.exists("test")) {
  dir.create("test")
}

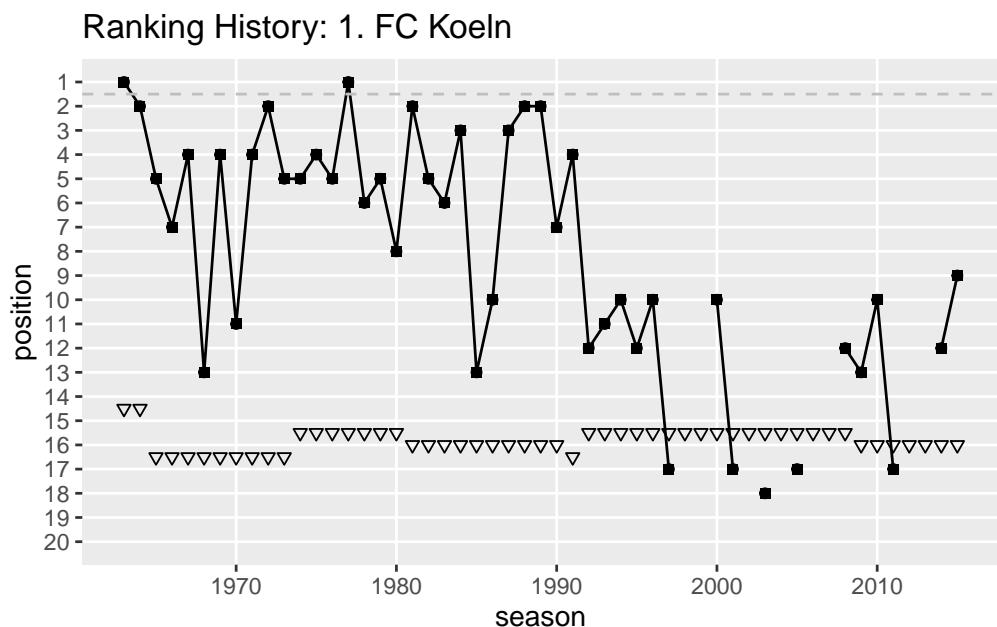
plots <- list()
for (club in unique(dfb_bal$team)) {
  dfb_subset <- subset(dfb_bal, team == club)

  p <- ggplot(dfb_subset, aes(x = season)) +
    geom_point(aes(y = position), shape = 15) +
    geom_line(aes(y = position)) +
    geom_point(aes(y = godown), shape = 25) +
    scale_y_reverse(breaks = seq(1, 20, by = 1), limits = c(20, 1)) +
    xlim(1963, 2015) +
    theme(panel.grid.minor = element_blank()) +
    geom_hline(yintercept = 1.5, linetype = "dashed", color = "gray") +
    geom_point(aes(y = position), shape = 1) +
    labs(title = paste("Ranking History:", club))
  ggsave(filename = paste("test/r_", club, ".png", sep = ""))
  plots[[club]] <- p
}

print(plots$`Meidericher SV`)
```



```
print(plots$`1. FC Koeln`)
```



```
# unload packages
suppressMessages(pacman::p_unload(
  bundesligR,
  tidyverse
))

# Remove the "test" directory and its contents after saving all graphs
unlink("test", recursive = TRUE)
```

## 9.16. Okun's Law

Suppose you aim to empirically examine unemployment and GDP for Germany and France. The data set that we use in the following is ‘forest.Rdata’ and should already been known to you from the lecture.

- (0) Write down your name, matriculation number, and date.
- (1) Set your working directory.
- (2) Clear your global environment.
- (3) Install and load the following packages: ‘tidyverse’, ‘sjPlot’, and ‘ggpubr’
- (4) Download and load the data, respectively, with the following code:

```
load(url("https://github.com/hubchev/courses/raw/main/dta/forest.Rdata"))
```

If that is not working, you can also download the data from ILIAS, save it in your working directory and load it from there with:

```
load("forest.Rdata")
```

- (5) Show the **first eight** observations of the dataset **df**.
- (6) Show the **last observation** of the dataset **df**.
- (7) Which type of data do we have here (Panel, cross-section, time series, ...)? Name the variable(s) that are necessary to identify the observations in the dataset.
- (8) Explain what the **assignment operator** in R is and what it is good for.
- (9) Write down the R code to store the number of observations and the number of variables that are in the dataset **df**. Name the object in which you store these numbers **observations\_df**.
- (10) In the dataset **df**, rename the variable ‘country.x’ to ‘nation’ and the variable ‘date’ to ‘year’.
- (11) Explain what the **pipe operator** in R is and what it is good for.
- (12) For the upcoming analysis you are only interested the following **variables** that are part of the dataframe **df**: nation, year, gdp, pop, gdppc, and unemployment. Drop all other variables from the dataframe **df**.
- (13) Create a variable that indicates the GDP per capita (‘gdp’ divided by ‘pop’). Name the variable ‘gdp\_pc’. (Hint: If you fail here, use the variable ‘gdppc’ which is already in the dataset as a replacement for ‘gdp\_pc’ in the following tasks.)
- (14) For the upcoming analysis you are only interested the following **countries** that are part of the dataframe **df**: Germany and France. Drop all other countries from the dataframe **df**.
- (15) Create a table showing the **average** unemployment rate and GDP per capita for Germany and France in the given years. Use the pipe operator. (Hint: See below for how your results should look like.)

## 9. Collection of exercises

```
# A tibble: 2 x 3
  nation `mean(unemployment)` `mean(gdppc)`
  <chr>          <dbl>        <dbl>
1 France         9.75       34356.
2 Germany        7.22       36739.
```

- (16) Create a table showing the unemployment rate and GDP per capita for Germany and France in the **year 2020**. Use the pipe operator. (Hint: See below for how your results should look like.)

```
# A tibble: 2 x 3
  nation `mean(unemployment)` `mean(gdppc)`
  <chr>          <dbl>        <dbl>
1 France         8.01       35786.
2 Germany        3.81       41315.
```

- (17) Create a table showing the **highest** unemployment rate and the **highest** GDP per capita for Germany and France during the given period. Use the pipe operator. (Hint: See below for how your results should look like.)

```
# A tibble: 2 x 3
  nation `max(unemployment)` `max(gdppc)`
  <chr>          <dbl>        <dbl>
1 France         12.6       38912.
2 Germany        11.2       43329.
```

- (18) Calculate the standard deviation of the unemployment rate and GDP per capita for Germany and France in the given years. (Hint: See below for how your result should look like.)

```
# A tibble: 2 x 3
  nation `sd(gdppc)` `sd(unemployment)`
  <chr>      <dbl>        <dbl>
1 France     2940.        1.58
2 Germany    4015.        2.37
```

- (19) In statistics, the coefficient of variation (COV) is a standardized measure of dispersion. It is defined as the ratio of the standard deviation ( $\sigma$ ) to the mean ( $\mu$ ):  $COV = \frac{\sigma}{\mu}$ . Write down the R code to calculate the coefficient of variation (COV) for the **unemployment rate** in Germany and France. (Hint: See below for what your result should look like.)

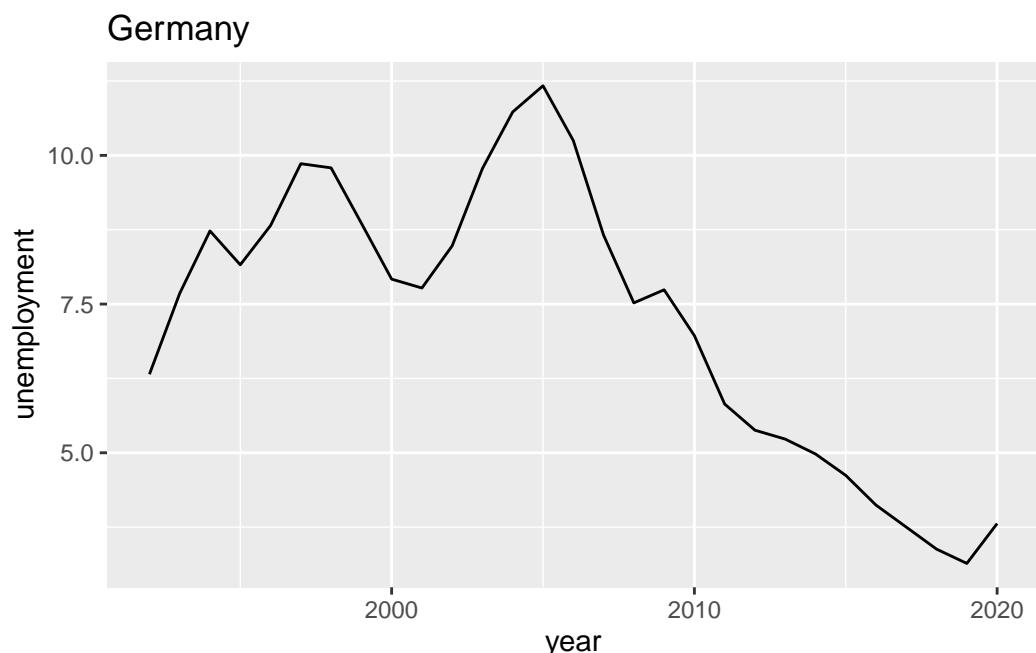
```
# A tibble: 2 x 4
  nation `sd(unemployment)` `mean(unemployment)` cov
  <chr>      <dbl>        <dbl> <dbl> <dbl>
1 France     1.58          9.75  0.162
2 Germany    2.37          7.22  0.328
```

- (20) Write down the R code to calculate the coefficient of variation (COV) for the **GDP per capita** in Germany and France. (Hint: See below for what your result should look like.)

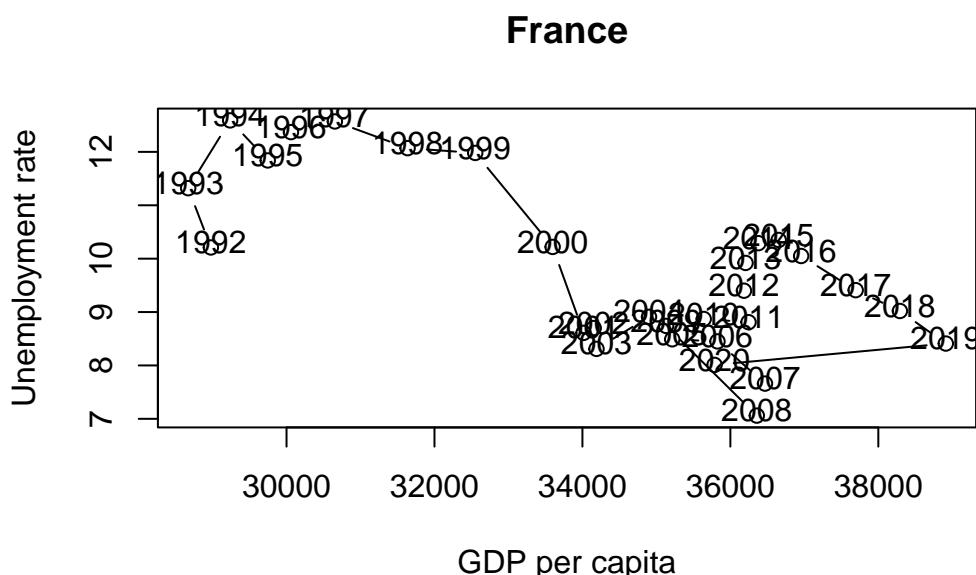
look like.)

```
# A tibble: 2 x 4
  nation `sd(gdppc)` `mean(gdppc)` cov
  <chr>     <dbl>        <dbl>   <dbl>
1 France      2940.       34356.  0.0856
2 Germany     4015.       36739.  0.109
```

- (21) Create a chart (bar chart, line chart, or scatter plot) that shows the unemployment rate of **Germany** over the available years. Label the chart ‘Germany’ with `ggtitle("Germany")`. Please note that you may choose any type of graphical representation. (Hint: Below you can see one of many possible examples of what your result may look like).



- (22) and 23. (*This task is worth 10 points*) The following chart shows the simultaneous development of the unemployment rate and GDP per capita over time for **France**.



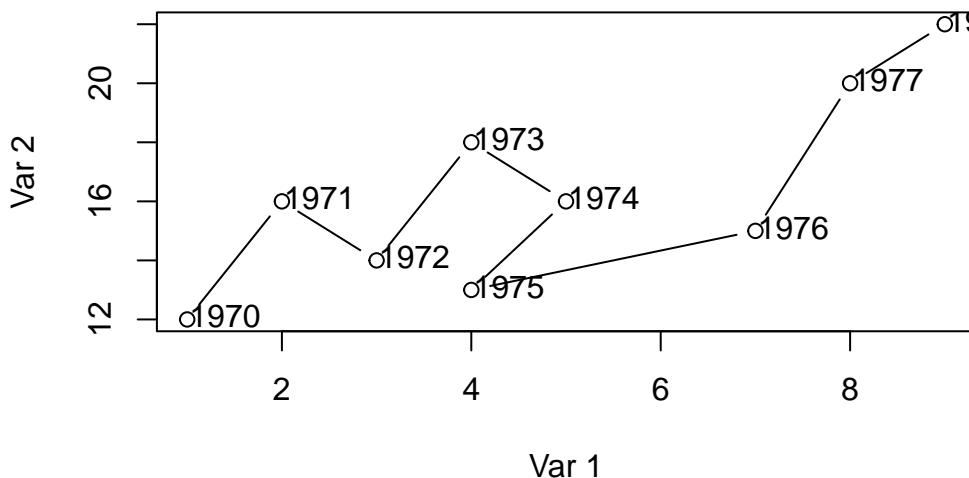
## 9. Collection of exercises

Suppose you want to visualize the simultaneous evolution of the unemployment rate and GDP per capita over time for Germany as well.

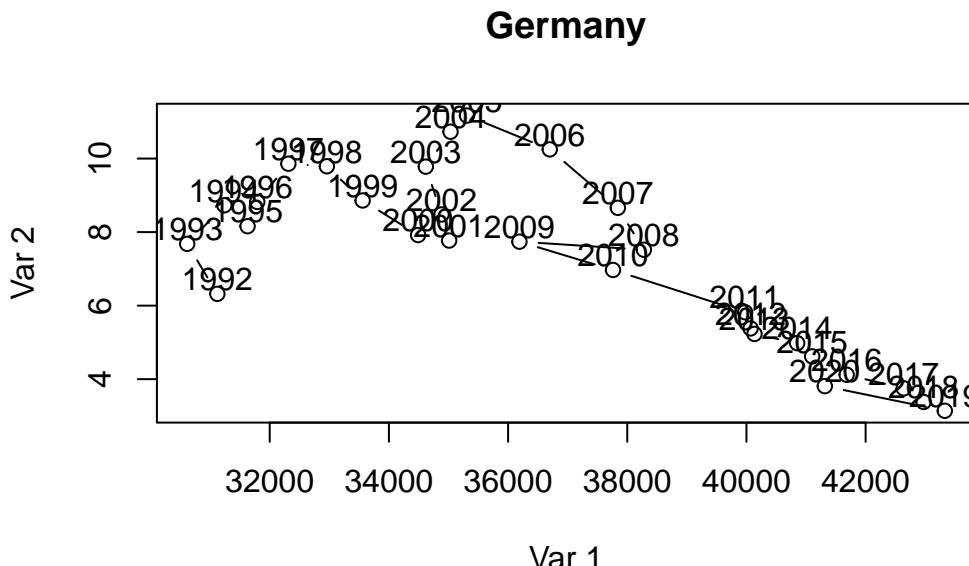
Suppose further that you have found the following lines of code that create the kind of chart you are looking for.

```
# Data
x <- c(1, 2, 3, 4, 5, 4, 7, 8, 9)
y <- c(12, 16, 14, 18, 16, 13, 15, 20, 22)
labels <- 1970:1978

# Connected scatter plot with text
plot(x, y, type = "b", xlab = "Var 1", ylab = "Var 2")
text(x + 0.4, y + 0.1, labels)
```



Use these lines of code and customize them to create the co-movement visualization for **Germany** using the available df data. The result should look something like this:



- (24) Interpret the two graphs above, which show the simultaneous evolution of the unemployment rate and GDP per capita over time for Germany and France. What are your expectations regarding the correlation between the unemployment rate and GDP per capita variables? Can you see this expectation in the figures? Discuss.

 Solution

The script uses the following functions: `aes`, `c`, `dim`, `filter`, `geom_line`, `ggplot`, `ggttitle`, `group_by`, `head`, `load`, `max`, `mean`, `mutate`, `plot`, `rename`, `sd`, `select`, `summarise`, `tail`, `text`, `title`, `url`.

**R script**

## 9. Collection of exercises

```
# setwd("/home/sthu/Dropbox/hsf/exams/22-11/scr/")

rm(list = ls())

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, ggpunr, sjPlot)

load(url("https://github.com/hubchev/courses/raw/main/dta/forest.Rdata"))

head(df, 8)

tail(df, 1)

# panel data set
# date and country.x

observations_df <- dim(df)

df <- rename(df, nation = country.x)
df <- rename(df, year = date)

df <- df |>
  select(nation, year, gdp, pop, gdppc, unemployment)

df <- df |>
  mutate(gdp_pc = gdp / pop)

df <- df |> filter(nation == "Germany" | nation == "France")

df |>
  group_by(nation) |>
  summarise(mean(unemployment), mean(gdppc))

df |>
  filter(year == 2020) |>
  group_by(nation) |>
  summarise(mean(unemployment), mean(gdppc))

df |>
  group_by(nation) |>
  summarise(max(unemployment), max(gdppc))

df |>
  group_by(nation) |>
  summarise(sd(gdppc), sd(unemployment))

df |>
  group_by(nation) |>
  summarise(sd(unemployment), mean(unemployment), cov = sd(unemployment) / mean(unemployment))

df |>
  group_by(nation) |>
  summarise(sd(gdppc), mean(gdppc), cov = sd(gdppc) / mean(gdppc))
224
df_pger <- df |>
```

## Output of the R script

```
# setwd("/home/sthu/Dropbox/hsf/exams/22-11/scr/")

rm(list = ls())

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, ggpunr, sjPlot)

load(url("https://github.com/hubchev/courses/raw/main/dta/forest.Rdata"))

head(df, 8)

# A tibble: 8 x 11
# Groups:   country.x [1]
  country.x     date     gdp gdp_growth unemployment region income forest    pop
  <chr>        <dbl>   <dbl>      <dbl>       <dbl> <chr>  <chr>   <dbl>   <dbl>
1 United Arab~ 1992 1.26e11     -2.48      1.84 Middl~ High ~  3.63 2.05e6
2 United Arab~ 1993 1.27e11     -4.34      1.85 Middl~ High ~  3.72 2.17e6
3 United Arab~ 1994 1.36e11      1.25      1.81 Middl~ High ~  3.81 2.29e6
4 United Arab~ 1995 1.45e11      1.35      1.80 Middl~ High ~  3.90 2.42e6
5 United Arab~ 1996 1.54e11      0.631     1.90 Middl~ High ~  3.99 2.54e6
6 United Arab~ 1997 1.66e11      2.83      1.98 Middl~ High ~  4.08 2.67e6
7 United Arab~ 1998 1.67e11     -4.77      2.14 Middl~ High ~  4.18 2.81e6
8 United Arab~ 1999 1.72e11     -2.40      2.22 Middl~ High ~  4.27 2.97e6
# i 2 more variables: unemployment_dif <dbl>, gdppc <dbl>

tail(df, 1)

# A tibble: 1 x 11
# Groups:   country.x [1]
  country.x     date     gdp gdp_growth unemployment region income forest    pop
  <chr>        <dbl>   <dbl>      <dbl>       <dbl> <chr>  <chr>   <dbl>   <dbl>
1 Zimbabwe    2020 1.94e10     -7.62      5.35 Sub-S~ Lower~  45.1 1.49e7
# i 2 more variables: unemployment_dif <dbl>, gdppc <dbl>
```

## 9. Collection of exercises

```
# panel data set
# date and country.x

observations_df <- dim(df)

df <- rename(df, nation = country.x)
df <- rename(df, year = date)

df <- df |>
  select(nation, year, gdp, pop, gdppc, unemployment)

df <- df |>
  mutate(gdp_pc = gdp / pop)

df <- df |> filter(nation == "Germany" | nation == "France")

df |>
  group_by(nation) |>
  summarise(mean(unemployment), mean(gdppc))

# A tibble: 2 x 3
  nation `mean(unemployment)` `mean(gdppc)`
  <chr>          <dbl>        <dbl>
1 France           9.75       34356.
2 Germany          7.22       36739.

df |>
  filter(year == 2020) |>
  group_by(nation) |>
  summarise(mean(unemployment), mean(gdppc))

# A tibble: 2 x 3
  nation `mean(unemployment)` `mean(gdppc)`
  <chr>          <dbl>        <dbl>
1 France           8.01       35786.
2 Germany          3.81       41315.

df |>
  group_by(nation) |>
  summarise(max(unemployment), max(gdppc))

# A tibble: 2 x 3
  nation `max(unemployment)` `max(gdppc)`
  <chr>          <dbl>        <dbl>
1 France           12.6       38912.
2 Germany          11.2       43329.

df |>
  group_by(nation) |>
  summarise(sd(gdppc), sd(unemployment))
```

## 9. Collection of exercises

```
# A tibble: 2 x 3
  nation `sd(gdppc)` `sd(unemployment)`
  <chr>      <dbl>          <dbl>
1 France       2940.         1.58
2 Germany      4015.        2.37

df |>
  group_by(nation) |>
  summarise(sd(unemployment), mean(unemployment), cov = sd(unemployment) / mean(unemployment))

# A tibble: 2 x 4
  nation `sd(unemployment)` `mean(unemployment)` cov
  <chr>      <dbl>           <dbl>    <dbl> <dbl>
1 France       1.58            9.75   0.162
2 Germany      2.37            7.22   0.328

df |>
  group_by(nation) |>
  summarise(sd(gdppc), mean(gdppc), cov = sd(gdppc) / mean(gdppc))

# A tibble: 2 x 4
  nation `sd(gdppc)` `mean(gdppc)` cov
  <chr>      <dbl>           <dbl>    <dbl>
1 France       2940.          34356.  0.0856
2 Germany      4015.          36739.  0.109

df_pger <- df |>
  filter(nation == "Germany")

pger <- ggplot(df_pger, aes(x = year, y = unemployment)) +
  geom_line() +
  ggtitle("Germany")

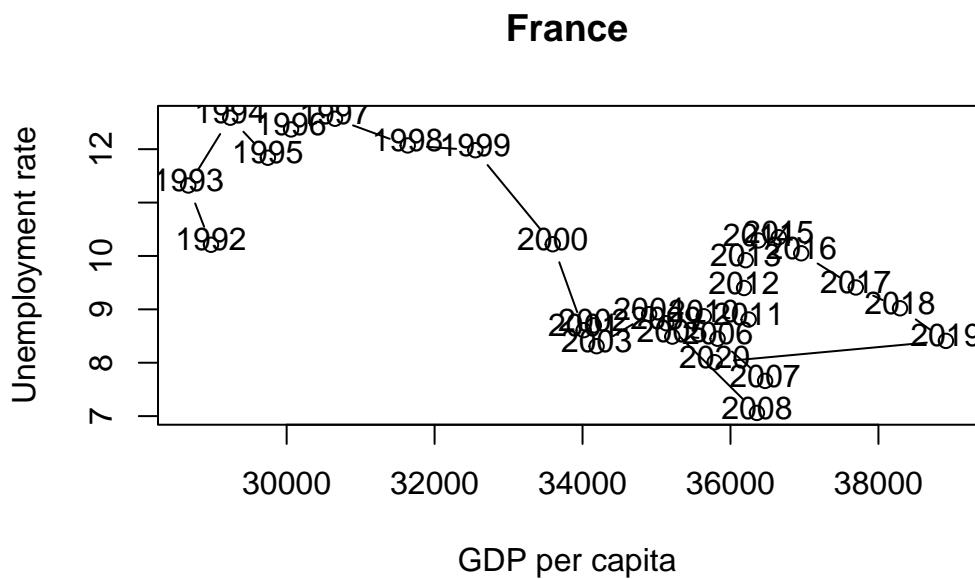
plot(pger)
```



```

labels <- 1992:2020
dfra <- df |> filter(nation == "France")
plot(dfra$gdppc, dfra$unemployment,
      type = "b",
      xlab = "GDP per capita", ylab = "Unemployment rate"
)
text(dfra$gdppc + 0.1, dfra$unemployment + 0.1, labels)
title("France")

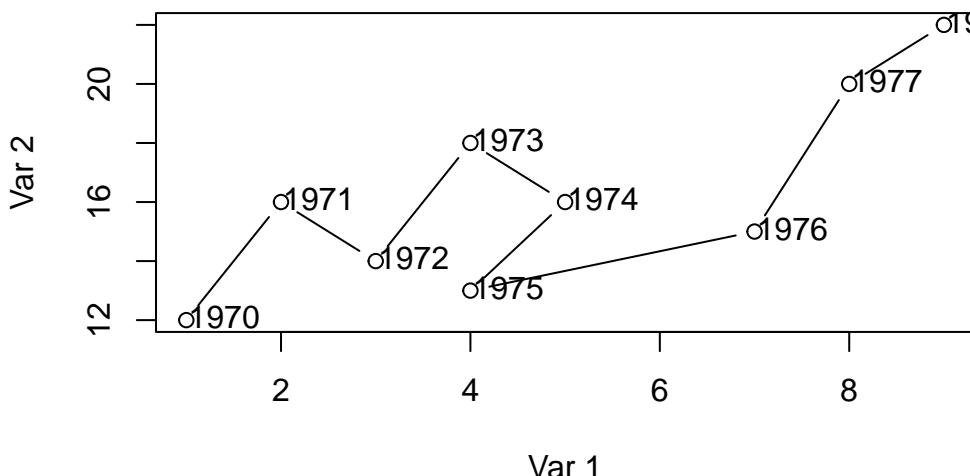
```



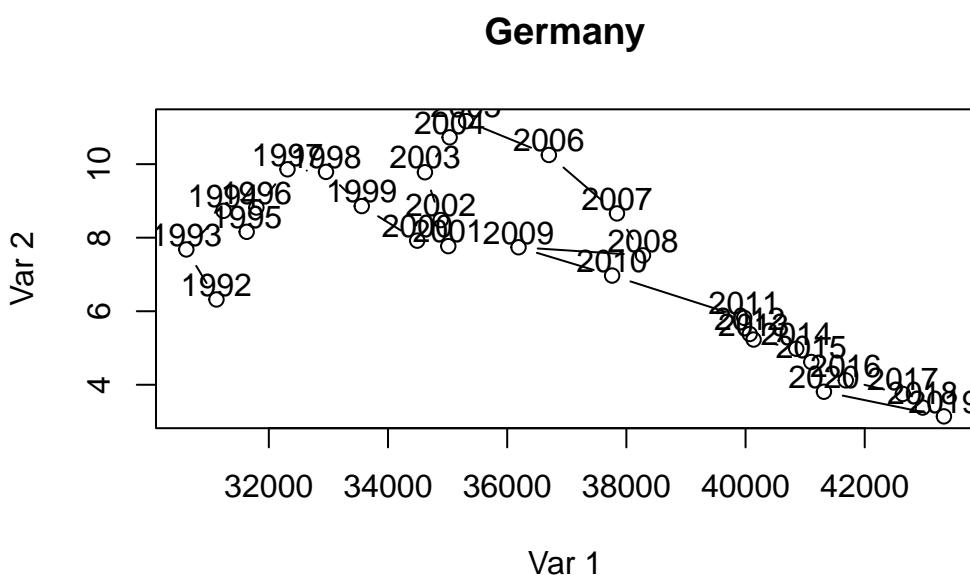
## 9. Collection of exercises

```
# Data
x <- c(1, 2, 3, 4, 5, 4, 7, 8, 9)
y <- c(12, 16, 14, 18, 16, 13, 15, 20, 22)
labels <- 1970:1978

# Connected scatter plot with text
plot(x, y, type = "b", xlab = "Var 1", ylab = "Var 2")
text(x + 0.4, y + 0.1, labels)
```



```
dfger <- df |> filter(nation == "Germany")
labels <- 1992:2020
plot(dfger$gdppc, dfger$unemployment,
      type = "b",
      xlab = "Var 1", ylab = "Var 2"
    )
text(dfger$gdppc + 0.7, dfger$unemployment + 0.4, labels)
title("Germany")
```



```
suppressMessages(pacman::p_unload(tidyverse, ggpublisher, sjPlot))
```

## 9.17. Names and duplicates

1. Load the required packages (`pacman`, `tidyverse`, `janitor`, `babynames`, `stringr`).
2. Load the dataset from the URL: [https://github.com/hubchev/courses/raw/main/dta/df\\_names.RData](https://github.com/hubchev/courses/raw/main/dta/df_names.RData). Make yourself familiar with the data.
3. After loading the dataset, remove all objects except `df_2022` and `df_2022_error`.
4. Reorder the data using the `relocate` function so that `surname`, `name`, and `age` appear first. Save the changed data in a tibble called `df`.
5. Sort the data according to `surname`, `name`, and `age`.
6. Make a variable named `born` that contains the year of birth. How is the `born` variable calculated?
7. Create a new variable named `id` that identifies each person by `surname`, `name`, and their birth year (`born`). Why is this identifier useful?
8. Investigate how the data is identified. Are there any duplicates? If so, can you think of strategies to identify and how to deal with these duplicates.
9. Unload the packages used in the script. Why is unloading packages considered good practice?

### Solution

The script uses the following functions: `anti_join`, `arrange`, `c`, `cur_group_id`, `desc`, `dim`, `distinct`, `filter`, `get_dupes`, `glimpse`, `group_by`, `head`, `load`, `max`, `mutate`, `n`, `paste`, `relocate`, `row_number`, `setdiff`, `summary`, `tail`, `ungroup`, `url`.

### R script

## 9. Collection of exercises

```
# Find duplicates

# set working directory
# setwd("~/Dropbox/hsf/test/initial_script")

# clear environment
rm(list = ls())

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, janitor, babynames, stringr)

load(url("https://github.com/hubchev/courses/raw/main/dta/df_names.RData"))

# Remove all objects except df_2022 and df_2022_error
rm(list = setdiff(ls(), c("df_2022_error", "df_2022")))

# Re-order the data so that surname, name, and age appears first.
# Save the changed data in a tibble called `df`.
df <- df_2022 |>
  relocate(surname, name, age)

# Sort the data according to surname, name, and age.
df <- df |>
  arrange(surname, name, age)

# Inspect df_2022 and df_2022_error
df
dim(df)
head(df)
tail(df)
glimpse(df)
summary(df)

df_2022_error

# Make a variable that contains the year of birth. Name the variable `born`#
# and new dataframe `df`.
df <- df_2022 |>
  mutate(born = time - age)

# Make a new variable that identifies each person by surname, name,
# and their birth born. Name the variable `id`.
df <- df |>
  mutate(id = paste(surname, name, born, sep = "_"))

# How many different groups do exist?
df <- df |>
  group_by(id) |>
  mutate(id_num = cur_group_id()) |>
  ungroup()

max(df$id_num)

# Show groups that exist more than once.
df <- df |>
```

## Output of the R script

```
# Find duplicates

# set working directory
# setwd("~/Dropbox/hsf/test/initial_script")

# clear environment
rm(list = ls())

# load packages
if (!require(pacman)) install.packages("pacman")
pacman::p_load(tidyverse, janitor, babynames, stringr)

load(url("https://github.com/hubchev/courses/raw/main/dta/df_names.RData"))

# Remove all objects except df_2022 and df_2022_error
rm(list = setdiff(ls(), c("df_2022_error", "df_2022")))

# Re-order the data so that surname, name, and age appears first.
# Save the changed data in a tibble called `df`.
df <- df_2022 |>
  relocate(surname, name, age)

# Sort the data according to surname, name, and age.
df <- df |>
  arrange(surname, name, age)

# Inspect df_2022 and df_2022_error
df

# A tibble: 1,018 x 8
  surname name      age sex      cm  time error error_desc
  <chr>   <chr>    <dbl> <chr> <dbl> <dbl> <dbl> <chr>
1 Adams    Adonnis    30 M    192  2022     0 <NA>
2 Adams    Adonnis    30 M    192  2022     1 duplicate
3 Adams    Aila       79 F    157  2022     0 <NA>
4 Adams    Avenelle   69 F    157  2022     0 <NA>
5 Adams    Brysan     39 M    192  2022     0 <NA>
6 Adams    Eona       84 F    157  2022     0 <NA>
7 Adams    Eveline    42 F    157  2022     0 <NA>
8 Adams    Faithe     17 F    172. 2022     0 <NA>
9 Adams    Ineisha    47 F    157  2022     0 <NA>
10 Adams   Kloeigh    31 F    157  2022     0 <NA>
# i 1,008 more rows

dim(df)
```

```
[1] 1018     8
```

## 9. Collection of exercises

```
head(df)
```

```
# A tibble: 6 x 8
  surname name      age sex      cm  time error error_desc
  <chr>   <chr>    <dbl> <chr> <dbl> <dbl> <dbl> <chr>
1 Adams   Adonnis    30 M     192  2022    0 <NA>
2 Adams   Adonnis    30 M     192  2022    1 duplicate
3 Adams   Aila       79 F     157  2022    0 <NA>
4 Adams   Avenelle   69 F     157  2022    0 <NA>
5 Adams   Brysan     39 M     192  2022    0 <NA>
6 Adams   Eona       84 F     157  2022    0 <NA>
```

```
tail(df)
```

```
# A tibble: 6 x 8
  surname name      age sex      cm  time error error_desc
  <chr>   <chr>    <dbl> <chr> <dbl> <dbl> <dbl> <chr>
1 Young   Leiliana   54 F     157  2022    0 <NA>
2 Young   Shamar     23 M     192  2022    0 <NA>
3 Young   Tajanay    1 F      81.5 2022    0 <NA>
4 huber   Stephan    186 M    41   2022    1 age/cm false, not capitalized ~
5 huber   Stephan    NA <NA>  NA   2022    1 wrong name
6 <NA>    Zita       6 <NA>   110  2022    2 surname missing, sex unspecifi~
```

```
glimpse(df)
```

Rows: 1,018

Columns: 8

```
$ surname    <chr> "Adams", "Adams", "Adams", "Adams", "Adams", "Adams", "Adam~
$ name       <chr> "Adonnis", "Adonnis", "Aila", "Avenelle", "Brysan", "Eona", ~
$ age        <dbl> 30, 30, 79, 69, 39, 84, 42, 17, 47, 31, 65, 80, 6, 5, 5, 20~
$ sex        <chr> "M", "M", "F", "F", "M", "F", "F", "F", "F", "M", "F", ~
$ cm         <dbl> 192.00000, 192.00000, 157.00000, 157.00000, 192.00000, 157.~
$ time       <dbl> 2022, 2022, 2022, 2022, 2022, 2022, 2022, 2022, ~
$ error      <dbl> 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, ~
$ error_desc <chr> NA, "duplicate", NA, NA~
```

```
summary(df)
```

surname	name	age	sex
Length:1018	Length:1018	Min. : 1.00	Length:1018
Class :character	Class :character	1st Qu.: 21.00	Class :character
Mode :character	Mode :character	Median : 43.00	Mode :character
		Mean : 45.75	
		3rd Qu.: 69.00	
		Max. : 399.00	
		NA's : 2	
cm	time	error	error_desc

## 9. Collection of exercises

```

Min.   : 41.0   Min.   :2022   Min.   :0.00000  Length:1018
1st Qu.:157.0  1st Qu.:2022  1st Qu.:0.00000  Class  :character
Median :157.0  Median :2022  Median :0.00000  Mode   :character
Mean   :163.2  Mean   :2022  Mean   :0.02456
3rd Qu.:192.0  3rd Qu.:2022 3rd Qu.:0.00000
Max.   :295.0  Max.   :2022  Max.   :3.00000
NA's    :4

```

### df\_2022\_error

```

# A tibble: 18 x 8
  sex     name    surname   age     cm  time error error_desc
  <chr>  <chr>   <chr>     <dbl> <dbl> <dbl> <dbl> <chr>
1 M      Savier  Campbell  72    192  2022    1 duplicate
2 F      Tina    Adams     5     98.0 2022    1 duplicate
3 F      Aberry  Allen    79    157  2022    1 duplicate
4 M      Adonnis Adams    30    192  2022    1 duplicate
5 M      Stephan Maier    41    186  2022    1 wrong surname
6 <NA>   Stephan huber  NA     NA   2022    1 wrong name
7 M      stephan Huber   186   41   2022    1 age/cm false, not capitalized-
8 M      Stephan huber   186   41   2022    1 age/cm false, not capitalized-
9 M      Stephan Huber   41    186  2022    1 duplicate
10 M     Stephan Huber   41    NA   2022    1 duplicate, cm NA
11 F     Rosa    Huber    9     NA   2022    3 only age and sex given
12 <NA>   Rosa    Huber   NA    130  2022    3 age missing, sex unspecified
13 <NA>   Ignaz   Huber   7     NA   2022    2 cm missing, sex unspecified
14 <NA>   Zita    <NA>     6    110  2022    2 surname missing, sex unspecified
15 <NA>   Alois   Huber   3    295  2022    2 cm not possible, sex unspecified
16 F      Martina Huber  399   169  2022    2 age not possible
17 M      Stephan Huber  41    186  2022    0 no error
18 M      Stephan Huber  41    186  2022    1 duplicate

# Make a variable that contains the year of birth. Name the variable `born`#
# and new dataframe `df` .
df <- df_2022 |>
  mutate(born = time - age)

# Make a new variable that identifies each person by surname, name,
# and their birth born. Name the variable `id` .
df <- df |>
  mutate(id = paste(surname, name, born, sep = "_"))

# How many different groups do exist?
df <- df |>
  group_by(id) |>
  mutate(id_num = cur_group_id()) |>
  ungroup()

max(df$id_num)

```

## 9. Collection of exercises

```
[1] 1011
```

```
# Show groups that exist more than once.
df <- df |>
  group_by(id) |>
  mutate(
    dup_count = row_number(),
    dup_sum   = n()
  ) |>
  ungroup() |>
  arrange(id)

df |> filter(dup_sum > 1)
```

```
# A tibble: 12 x 13
  sex   name   surname   age     cm   time error error_desc   born id   id_num
  <chr> <chr> <chr>   <dbl> <dbl> <dbl> <dbl> <chr>   <dbl> <chr> <int>
1 M     Adonnis Adams     30 192   2022     0 <NA>   1992 Adam~   1
2 M     Adonnis Adams     30 192   2022     1 duplicate 1992 Adam~   1
3 F     Tina   Adams      5  98.0  2022     1 duplicate 2017 Adam~   13
4 F     Tina   Adams      5  98.0  2022     0 <NA>   2017 Adam~   13
5 F     Abery  Allen     79 157   2022     0 <NA>   1943 Alle~   15
6 F     Abery  Allen     79 157   2022     1 duplicate 1943 Alle~   15
7 M     Savier Campbell   72 192   2022     0 <NA>   1950 Camp~  100
8 M     Savier Campbell   72 192   2022     1 duplicate 1950 Camp~  100
9 M     Stephan Huber    41 186   2022     1 duplicate 1981 Hube~  383
10 M    Stephan Huber    41 186   2022     0 no error 1981 Hube~  383
11 M    Stephan Huber    41 186   2022     1 duplicate 1981 Hube~  383
12 M    Stephan Huber    41  NA    2022     1 duplicate,~ 1981 Hube~  383
# i 2 more variables: dup_count <int>, dup_sum <int>
```

```
df |> get_dupes(name, surname)
```

```
# A tibble: 18 x 14
  name   surname dupe_count sex     age     cm   time error error_desc   born id
  <chr> <chr>       <int> <chr> <dbl> <dbl> <dbl> <chr>   <dbl> <chr>
1 Step~ Huber        4 M     41 186   2022     1 duplicate 1981 Hube~
2 Step~ Huber        4 M     41 186   2022     0 no error 1981 Hube~
3 Step~ Huber        4 M     41 186   2022     1 duplicate 1981 Hube~
4 Step~ Huber        4 M     41  NA    2022     1 duplicate~ 1981 Hube~
5 Abery Allen        2 F     79 157   2022     0 <NA>   1943 Alle~
6 Abery Allen        2 F     79 157   2022     1 duplicate 1943 Alle~
7 Adon~ Adams         2 M     30 192   2022     0 <NA>   1992 Adam~
8 Adon~ Adams         2 M     30 192   2022     1 duplicate 1992 Adam~
9 Merl~ Miller        2 F     12 153.  2022     0 <NA>   2010 Mill~
10 Merl~ Miller       2 F     2  99.9  2022     0 <NA>   2020 Mill~
11 Rosa  Huber        2 F     9  NA    2022     3 only age ~ 2013 Hube~
12 Rosa  Huber        2 <NA>  NA 130    2022     3 age missi~  NA Hube~
13 Savi~ Campbe~       2 M     72 192   2022     0 <NA>   1950 Camp~
```

## 9. Collection of exercises

```

14 Savi~ Campbe~      2 M      72 192    2022      1 duplicate   1950 Camp~
15 Step~ huber        2 M      186 41     2022      1 age/cm fa~  1836 huber~
16 Step~ huber        2 <NA>    NA  NA     2022      1 wrong name  NA huber~
17 Tina Adams         2 F       5  98.0    2022      1 duplicate   2017 Adam~
18 Tina Adams         2 F       5  98.0    2022      0 <NA>       2017 Adam~

# i 3 more variables: id_num <int>, dup_count <int>, dup_sum <int>

```

```
# Make yourself familiar with the function `get_dupes()` from `janitor` package.
df |> get_dupes()
```

No variable names specified - using all columns.

No duplicate combinations found of: sex, name, surname, age, cm, time, error, error\_desc

```

# A tibble: 0 x 14
# i 14 variables: sex <chr>, name <chr>, surname <chr>, age <dbl>, cm <dbl>,
#   time <dbl>, error <dbl>, error_desc <chr>, born <dbl>, id <chr>,
#   id_num <int>, dup_count <int>, dup_sum <int>, dupe_count <int>
```

```
df |> get_dupes(surname, name)
```

# A tibble: 18 x 14

	surname	name	dupe_count	sex	age	cm	time	error	error_desc	born	id
	<chr>	<chr>	<int>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>	<dbl>	<chr>
1	Huber	Step~	4	M	41	186	2022	1	duplicate	1981	Huber~
2	Huber	Step~	4	M	41	186	2022	0	no error	1981	Huber~
3	Huber	Step~	4	M	41	186	2022	1	duplicate	1981	Huber~
4	Huber	Step~	4	M	41	NA	2022	1	duplicate~	1981	Huber~
5	Adams	Adon~	2	M	30	192	2022	0	<NA>	1992	Adam~
6	Adams	Adon~	2	M	30	192	2022	1	duplicate	1992	Adam~
7	Adams	Tina	2	F	5	98.0	2022	1	duplicate	2017	Adam~
8	Adams	Tina	2	F	5	98.0	2022	0	<NA>	2017	Adam~
9	Allen	Abery	2	F	79	157	2022	0	<NA>	1943	Alle~
10	Allen	Abery	2	F	79	157	2022	1	duplicate	1943	Alle~
11	Campbe~	Savi~	2	M	72	192	2022	0	<NA>	1950	Camp~
12	Campbe~	Savi~	2	M	72	192	2022	1	duplicate	1950	Camp~
13	Huber	Rosa	2	F	9	NA	2022	3	only age ~	2013	Huber~
14	Huber	Rosa	2	<NA>	NA	130	2022	3	age missi~	NA	Huber~
15	Miller	Merl~	2	F	12	153.	2022	0	<NA>	2010	Mill~
16	Miller	Merl~	2	F	2	99.9	2022	0	<NA>	2020	Mill~
17	huber	Step~	2	M	186	41	2022	1	age/cm fa~	1836	huber~
18	huber	Step~	2	<NA>	NA	NA	2022	1	wrong name	NA	huber~

```
# i 3 more variables: id_num <int>, dup_count <int>, dup_sum <int>
```

```
df |> get_dupes(id)
```

# A tibble: 12 x 14

id	dupe_count	sex	name	surname	age	cm	time	error	error_desc	born
----	------------	-----	------	---------	-----	----	------	-------	------------	------

## 9. Collection of exercises

```

<chr>     <int> <chr> <chr> <dbl> <dbl> <dbl> <dbl> <chr>     <dbl>
1 Hube~      4 M    Step~ Huber      41 186   2022     1 duplicate  1981
2 Hube~      4 M    Step~ Huber      41 186   2022     0 no error  1981
3 Hube~      4 M    Step~ Huber      41 186   2022     1 duplicate  1981
4 Hube~      4 M    Step~ Huber      41 NA    2022     1 duplicate~ 1981
5 Adam~      2 M    Adon~ Adams     30 192   2022     0 <NA>    1992
6 Adam~      2 M    Adon~ Adams     30 192   2022     1 duplicate  1992
7 Adam~      2 F    Tina Adams     5 98.0  2022     1 duplicate  2017
8 Adam~      2 F    Tina Adams     5 98.0  2022     0 <NA>    2017
9 Alle~       2 F    Abery Allen    79 157   2022     0 <NA>    1943
10 Alle~      2 F    Abery Allen    79 157   2022     1 duplicate  1943
11 Camp~      2 M    Savi~ Campbell 72 192   2022     0 <NA>    1950
12 Camp~      2 M    Savi~ Campbell 72 192   2022     1 duplicate  1950
# i 3 more variables: id_num <int>, dup_count <int>, dup_sum <int>
```

```

df_uni <- df |>
  arrange() |>
  distinct(id, .keep_all = TRUE)

df_uni_b <- df |>
  arrange(desc(dup_count)) |>
  distinct(id, .keep_all = TRUE)

anti_join(df, df_uni)
```

Joining with `by = join\_by(sex, name, surname, age, cm, time, error, error\_desc, born, id, id\_num, dup\_count, dup\_sum)`

```

# A tibble: 7 x 13
  sex   name   surname   age     cm   time error error_desc   born id   id_num
  <chr> <chr> <chr>   <dbl> <dbl> <dbl> <dbl> <chr>   <dbl> <chr> <int>
1 M    Adonnis Adams     30 192   2022     1 duplicate  1992 Adam~    1
2 F    Tina   Adams      5 98.0  2022     0 <NA>    2017 Adam~    13
3 F    Abery  Allen     79 157   2022     1 duplicate  1943 Alle~    15
4 M    Savier Campbell   72 192   2022     1 duplicate  1950 Camp~   100
5 M    Stephan Huber    41 186   2022     0 no error  1981 Hube~   383
6 M    Stephan Huber    41 186   2022     1 duplicate  1981 Hube~   383
7 M    Stephan Huber    41 NA    2022     1 duplicate, ~  1981 Hube~   383
# i 2 more variables: dup_count <int>, dup_sum <int>
```

```
anti_join(df, df_uni_b)
```

Joining with `by = join\_by(sex, name, surname, age, cm, time, error, error\_desc, born, id, id\_num, dup\_count, dup\_sum)`

```

# A tibble: 7 x 13
  sex   name   surname   age     cm   time error error_desc   born id   id_num
  <chr> <chr> <chr>   <dbl> <dbl> <dbl> <dbl> <chr>   <dbl> <chr> <int>
1 M    Adonnis Adams     30 192   2022     0 <NA>    1992 Adams~    1
```

```

2 F Tina Adams      5 98.0 2022    1 duplicate  2017 Adams_~ 13
3 F Abery Allen    79 157 2022     0 <NA>    1943 Allen_~ 15
4 M Savier Campbell 72 192 2022     0 <NA>    1950 Campbe~ 100
5 M Stephan Huber   41 186 2022    1 duplicate  1981 Huber_~ 383
6 M Stephan Huber   41 186 2022    0 no error  1981 Huber_~ 383
7 M Stephan Huber   41 186 2022    1 duplicate  1981 Huber_~ 383
# i 2 more variables: dup_count <int>, dup_sum <int>

# unload packages
suppressMessages(pacman::p_unload(tidyverse, janitor, babynames, stringr))

```

## 9.18. Zipf's law

The data under investigation includes population information for various German cities, identified by the variable `stadt`, spanning the years 1970, 1987, and 2010. The variable `status` provides details about the legislative status of the cities, and the variable `state` (Bundesland) indicates the state in which each respective city is situated.

### Preamble

- (1) Set your working directory.
- (2) Clear your global environment.
- (3) Install and load the following packages: ‘tidyverse’, ‘haven’, and ‘janitor’.

### Read in, inspect, and clean the data

- (4) Download and load the data, respectively, with the following code:

```

df <- read_dta(
  "https://github.com/hubchev/courses/raw/main/dta/city.dta",
  encoding = "latin1"
) |>
  as_tibble()

```

If that is not working, you can also download the data from ILIAS, save it in your working directory and load it from there with:

```
load("city.RData")
```

- (5) Show the first six and the last six observations of the dataset `df`.
- (6) How many observations (rows) and variables (columns) are in the dataset?
- (7) Show for all numerical variables the summary statistics including the mean, median, minimum, and the maximum.
- (8) Rename the variable `stadt` to `city`.
- (9) Remove the variables `pop1970` and `pop1987`.
- (10) Replicate the following table which contains some summary statistics.

## 9. Collection of exercises

```
# A tibble: 17 x 3
  state      `mean(pop2011)` `sum(pop2011)`
  <chr>          <dbl>        <dbl>
1 Baden-Wrttemberg     7580        7580
2 Baden-Württemberg   23680.      7837917
3 Bayern             23996.      7558677
4 Berlin              3292365    3292365
5 Brandenburg         18472.      1865632
6 Bremen              325432.     650863
7 Hamburg             1706696    1706696
8 Hessen              22996.      5036121
9 Mecklenburg-Vorpommern 27034.     811005
10 Niedersachsen      24107.      6219515
11 Nordrhein-Westfalen 47465.      18036727
12 Rheinland-Pfalz    25644.      1871995
13 Saarland            NA          NA
14 Sachsen             27788.      2973351
15 Sachsen-Anhalt     21212.      1993915
16 Schleswig-Holstein  24157.      1739269
17 Th_ringen           29192.     1167692
```

- (11) The states “Baden-Wrttemberg” and “Th\_ringen” are falsely pronounced. Correct the names and regenerate the summary statistics table presented above. Your result should look like this:

```
# A tibble: 16 x 3
  state      `mean(pop2011)` `sum(pop2011)`
  <chr>          <dbl>        <dbl>
1 Baden-Württemberg   23631.      7845497
2 Bayern             23996.      7558677
3 Berlin              3292365    3292365
4 Brandenburg         18472.      1865632
5 Bremen              325432.     650863
6 Hamburg             1706696    1706696
7 Hessen              22996.      5036121
8 Mecklenburg-Vorpommern 27034.     811005
9 Niedersachsen      24107.      6219515
10 Nordrhein-Westfalen 47465.      18036727
11 Rheinland-Pfalz    25644.      1871995
12 Saarland            NA          NA
13 Sachsen             27788.      2973351
14 Sachsen-Anhalt     21212.      1993915
15 Schleswig-Holstein  24157.      1739269
16 Thüringen           29192.     1167692
```

- (12) To investigate the reason for observing only NAs for Saarland, examine all cities within Saarland. Therefore, please display all observations for cities in Saarland in the Console, as illustrated below.

```
# A tibble: 47 x 5
  city      status  state  pop2011 rankX
```

9. Collection of exercises

<chr>	<chr>	<chr>	<dbl>	<dbl>
1 Perl	Commune	Saarland	7775	2003
2 Freisen	Commune	Saarland	8270	1894
3 Großrosseln	Commune	Saarland	8403	1868
4 Nonnweiler	Commune	Saarland	8844	1775
5 Nalbach	Commune	Saarland	9302	1678
6 Wallerfangen	Commune	Saarland	9542	1642
7 Kirkel	Commune	Saarland	10058	1541
8 Merchweiler	Commune	Saarland	10219	1515
9 Nohfelden	Commune	Saarland	10247	1511
10 Friedrichsthal	City	Saarland	10409	1489
11 Marpingen	Commune	Saarland	10590	1461
12 Mandelbachtal	Commune	Saarland	11107	1390
13 Kleinblittersdorf	Commune	Saarland	11396	1354
14 Überherrn	Commune	Saarland	11655	1317
15 Mettlach	Commune	Saarland	12180	1241
16 Tholey	Commune	Saarland	12385	1217
17 Saarwellingen	Commune	Saarland	13348	1104
18 Quierschied	Commune	Saarland	13506	1088
19 Spiesen-Elversberg	Commune	Saarland	13509	1086
20 Rehlingen-Siersburg	Commune	Saarland	14526	996
21 Riegelsberg	Commune	Saarland	14763	982
22 Ottweiler	City	Saarland	14934	969
23 Beckingen	Commune	Saarland	15355	931
24 Losheim am See	Commune	Saarland	15906	887
25 Schiffweiler	Commune	Saarland	15993	882
26 Wadern	City	Saarland	16181	874
27 Schmelz	Commune	Saarland	16435	857
28 Sulzbach/Saar	City	Saarland	16591	849
29 Illingen	Commune	Saarland	16978	827
30 Schwalbach	Commune	Saarland	17320	812
31 Eppelborn	Commune	Saarland	17726	793
32 Wadgassen	Commune	Saarland	17885	785
33 Bexbach	City	Saarland	18038	777
34 Heusweiler	Commune	Saarland	18201	762
35 Püttlingen	City	Saarland	19134	718
36 Lebach	City	Saarland	19484	701
37 Dillingen/Saar	City	Saarland	20253	654
38 Blieskastel	City	Saarland	21255	601
39 St. Wendel	City	Saarland	26220	460
40 Merzig	City	Saarland	29727	392
41 Saarlouis	City	Saarland	34479	323
42 St. Ingbert	City	Saarland	36645	299
43 Völklingen	City	Saarland	38809	279
44 Homburg	City	Saarland	41502	247
45 Neunkirchen	City	Saarland	46172	206
46 Saarbrücken	City	Saarland	175853	43
47 Perl	Commune	Saarland	NA	NA

- (13) With reference to the table above, we have identified an entry for the city of Perl that solely consists of NAs. This city is duplicated in the dataset, appearing at positions 1 and 47. The latter duplicate contains only NAs and can be safely removed without the loss

## 9. Collection of exercises

of valuable information. Please eliminate this duplication and regenerate the list of all cities in the Saarland.

- (14) Calculate the total population and average size of cities in Saarland.
- (15) Check if any other city is recorded more than once in the dataset. To do so, reproduce the table below.

# A tibble: 23 x 5	# Groups: city [11]	city	status	state	pop2011	unique_count
		<chr>	<chr>	<chr>	<dbl>	<int>
1	Bonn	Bonn	City with County Rights	Nordrhein-Westfalen	305765	3
2	Bonn	Bonn	City with County Rights	Nordrhein-Westfalen	305765	3
3	Bonn	Bonn	City with County Rights	Nordrhein-Westfalen	305765	3
4	Brühl	Brühl	Commune	Baden-Württemberg	13805	2
5	Brühl	Brühl	City	Nordrhein-Westfalen	43568	2
6	Erbach	Erbach	City	Baden-Württemberg	13024	2
7	Erbach	Erbach	City	Hessen	13245	2
8	Fürth	Fürth	City with County Rights	Bayern	115613	2
9	Fürth	Fürth	Commune	Hessen	10481	2
10	Lichtenau	Lichtenau	City	Nordrhein-Westfalen	10473	2
11	Lichtenau	Lichtenau	Commune	Sachsen	7544	2
12	Münster	Münster	Commune	Hessen	14071	2
13	Münster	Münster	City with County Rights	Nordrhein-Westfalen	289576	2
14	Neunkirchen	Neunkirchen	Commune	Nordrhein-Westfalen	13930	2
15	Neunkirchen	Neunkirchen	City	Saarland	46172	2
16	Neuried	Neuried	Commune	Baden-Württemberg	9383	2
17	Neuried	Neuried	Commune	Bayern	8277	2
18	Petersberg	Petersberg	Commune	Hessen	14766	2
19	Petersberg	Petersberg	Commune	Sachsen-Anhalt	10097	2
20	Senden	Senden	City	Bayern	21560	2
21	Senden	Senden	Commune	Nordrhein-Westfalen	19976	2
22	Staufenberg	Staufenberg	City	Hessen	8114	2
23	Staufenberg	Staufenberg	Commune	Niedersachsen	7983	2

- (16) The table indicates that the city of Bonn appears three times in the dataset, and all three observations contain identical information. Thus, remove two of these observations to ensure that Bonn is uniquely represented in the dataset. All other cities that occur more than once in the data are situated in different states. That means, these are distinct cities that coincidentally share the same name.

### Data analysis (Zipf's Law)

**\*Note:** If you have failed to solve the data cleaning tasks above, you can download the cleaned data from ILIAS, save it in your working directory and load it from there with:  
`load("city_clean.RData")`

In the following, you aim to examine the validity of Zipf's Law for Germany. Zipf's Law postulates how the size of cities is distributed. The "law" states that there is a special relationship between the size of a city and the rank it occupies in a series sorted by city size. In the estimation equation

$$\log(M_j) = c - q \log(R_j),$$

## 9. Collection of exercises

the law postulates a coefficient of ( $q = 1$ ).  $c$  is a constant;  $M_j$  is the size of city  $j$ ;  $R_j$  is the rank that city  $j$  occupies in a series sorted by city size.

(17)

Create a variable named `rank` that includes a ranking of cities based on the population size in the year 2011. Therefore, Berlin should have a rank of 1, Hamburg a rank of 2, Munich a rank of 3, and so on.

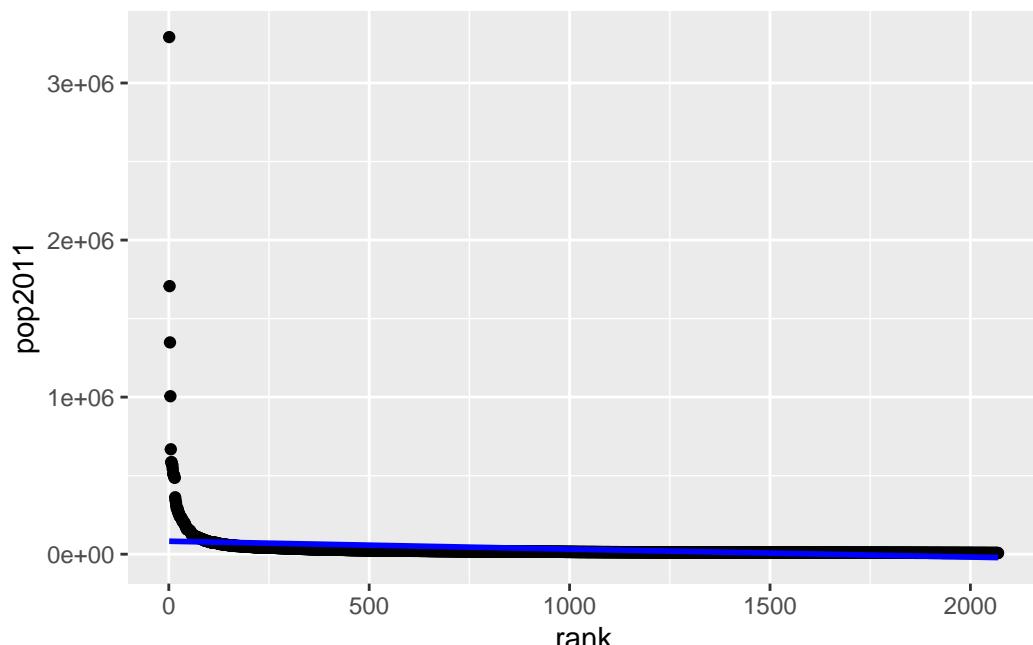
**Note:** If you cannot solve this task, use the variable `rankX` as a substitute for the variable `rank` that was not generated.

```
# A tibble: 6 x 3
  city          pop2011   rank
  <chr>        <dbl>     <int>
1 Berlin      3292365     1
2 Hamburg    1706696     2
3 München [Munich] 1348335     3
4 Köln [Cologne] 1005775     4
5 Frankfurt am Main 667925     5
6 Düsseldorf [Düsseldorf] 586291     6
```

(18) Calculate the Pearson Correlation Coefficient of the two variables `pop2011` and `rank`. The result should be:

```
[1] -0.2948903
```

(19) Create a scatter plot. On the x-axis, plot the variable `rank`, and on the y-axis, plot `pop2011`. Add a regression line representing the observed relationship to the same scatter plot.



(20) Logarithmize the variables `rank` and `pop2011`. Title the new variables as `lnrank` and `lnpop2011`, respectively. Here is a snapshot of the resulting variables:

## 9. Collection of exercises

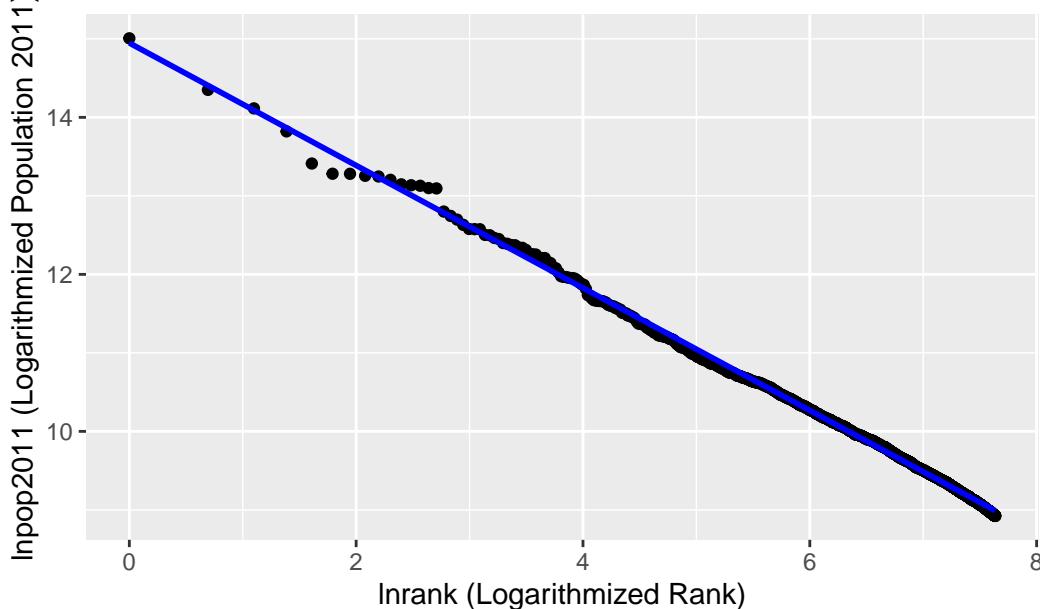
```
# A tibble: 6 x 5
  city              rank lnrank pop2011 lnpop2011
  <chr>           <int>   <dbl>    <dbl>     <dbl>
1 Berlin            1      0     3292365    15.0
2 Hamburg           2      0.693   1706696    14.4
3 München [Munich] 3      1.10    1348335    14.1
4 Köln [Cologne]   4      1.39    1005775    13.8
5 Frankfurt am Main 5      1.61    667925     13.4
6 Düsseldorf [Dusseldorf] 6      1.79    586291     13.3
```

- (21) Calculate the Pearson Correlation Coefficient of the two variables `lnpop2011` and `lnrank`.  
The result should be:

```
[1] -0.9990053
```

- (22) Create a scatter plot. On the x-axis, plot the variable `lnrank`, and on the y-axis, plot `lnpop2011`. Add a regression line representing the observed relationship to the same scatter plot. Additionally, add a title and label the axes like is shown here:

**Scatterplot with Regression Line**



- (23) Now, test the relationship postulated in Zipf's Law. Regress the logarithmic city size in the year 2011 on the logarithmic rank of a city in a series sorted by city size. Briefly interpret the results, addressing the coefficient of determination. Show the regression results. Here is one way to present the results of the regression (*Note: The way how you present your regression results do not matter*):

Call:  
`lm(formula = lnpop2011 ~ lnrank, data = df)`

Residuals:

Min	1Q	Median	3Q	Max
-0.28015	-0.01879	0.01083	0.02005	0.25973

## 9. Collection of exercises

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)    
(Intercept) 14.947859   0.005141   2908   <2e-16 ***
lnrank       -0.780259   0.000766  -1019   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.03454 on 2067 degrees of freedom
Multiple R-squared:  0.998, Adjusted R-squared:  0.998 
F-statistic: 1.038e+06 on 1 and 2067 DF,  p-value: < 2.2e-16
```

(24) Explain the following lines of code.

```
df <- df |>
  mutate(prediction = predict(zipf, newdata = df)) |>
  mutate(pred_pop = exp(prediction))
df |>
  select(city, pop2011, pred_pop) |>
  filter(city == "Regensburg")
```

```
# A tibble: 1 x 3
  city      pop2011 pred_pop
  <chr>     <dbl>    <dbl>
1 Regensburg 135403  134194.
```

### Solution

The script uses the following functions: `aes`, `arrange`, `as_tibble`, `c`, `case_when`, `cor`, `desc`, `dim`, `exp`, `filter`, `geom_point`, `geom_smooth`, `ggplot`, `group_by`, `head`, `is.na`, `labs`, `lm`, `log`, `mean`, `mutate`, `n`, `predict`, `print`, `read_dta`, `rename`, `row_number`, `save`, `select`, `starts_with`, `sum`, `summarise`, `summary`, `tail`, `ungroup`.

### R script

## 9. Collection of exercises

```
# load packages
if (!require(pacman)) install.packages("pacman")
suppressMessages(pacman::p_unload(all))
# setwd("~/Dropbox/hsf/exams/24-01/Rmd")

rm(list = ls())

pacman::p_load(tidyverse, haven, janitor, jtools)

df <- read_dta("https://github.com/hubchev/courses/raw/main/dta/city.dta",
  encoding = "latin1"
) |>
  as_tibble()

head(df)
tail(df)

dim(df)

summary(df)

df <- df |>
  rename(city = stadt)

df <- df |>
  select(-pop1970, -pop1987)

df |>
  group_by(state) |>
  summarise(
    mean(pop2011),
    sum(pop2011)
  )

df <- df |>
  mutate(state = case_when(
    state == "Baden-Wrttemberg" ~ "Baden-Württemberg",
    state == "Th_ringen" ~ "Thüringen",
    TRUE ~ state
  ))

df |>
  group_by(state) |>
  summarise(
    mean(pop2011),
    sum(pop2011)
  )

df |>
  filter(state == "Saarland") |>
  print(n = 100)

df <- df |>
  filter(!(city == "Perl" & is.na(pop2011)))
245
df |>
```

## Output of the R script

```
# load packages
if (!require(pacman)) install.packages("pacman")
suppressMessages(pacman::p_unload(all))
# setwd("~/Dropbox/hsf/exams/24-01/Rmd")

rm(list = ls())

pacman::p_load(tidyverse, haven, janitor, jtools)

df <- read_dta("https://github.com/hubchev/courses/raw/main/dta/city.dta",
  encoding = "latin1"
) |>
  as_tibble()

head(df)

# A tibble: 6 x 7
  stadt      status    state    pop1970  pop1987  pop2011 rankX
  <chr>     <chr>    <chr>    <dbl>    <dbl>    <dbl>    <dbl>
1 Vohenstrauß City     Bayern    7349     7059     7500    2069
2 Stockstadt a. Main Commune Bayern   6416     6615     7504    2068
3 Jesteburg     Commune Niedersachsen 4141     5818     7510    2067
4 Bordesholm    Commune Schleswig-Holstein 6011     6726     7513    2066
5 Herrieden     City     Bayern    5631     6250     7516    2065
6 Weida        City     Th_ringen NA       NA       7522    2064

tail(df)

# A tibble: 6 x 7
  stadt      status    state    pop1970  pop1987  pop2011 rankX
  <chr>     <chr>    <chr>    <dbl>    <dbl>    <dbl>    <dbl>
1 Frankfurt am Main City with County Rights Hessen 699297  618266  667925    5
2 Köln [Cologne]   City with County Rights Nordr~ 994705  928309 1005775    4
3 München [Munich] City with County Rights Bayern 1293599 1185421 1348335    3
4 Hamburg        City with County Rights Hambu~ 1793823 1592770 1706696    2
5 Berlin          City with County Rights Berlin 3210000 3260000 3292365    1
6 Perl            Commune           Saarl~      NA       NA       NA       NA

dim(df)

[1] 2072     7

summary(df)

  stadt      status    state    pop1970 
Length:2072  Length:2072  Length:2072  Min.   : 1604
```

## 9. Collection of exercises

```

Class :character   Class :character   Class :character   1st Qu.: 8149
Mode  :character   Mode  :character   Mode  :character   Median : 11912
                                         Mode  :character   Mean   : 30504
                                         Mode  :character   3rd Qu.: 21318
                                         Mode  :character   Max.   :3210000
                                         NA's   :355

```

pop1987	pop2011	rankX
Min. : 4003	Min. : 7500	Min. : 1.0
1st Qu.: 9194	1st Qu.: 9998	1st Qu.: 516.5
Median : 13118	Median : 13937	Median :1034.0
Mean   : 30854	Mean   : 30772	Mean   :1034.0
3rd Qu.: 23074	3rd Qu.: 24096	3rd Qu.:1551.5
Max.   :3260000	Max.   :3292365	Max.   :2069.0
NA's   :248	NA's   :1	NA's   :1

```

df <- df |>
  rename(city = stadt)

df <- df |>
  select(-pop1970, -pop1987)

df |>
  group_by(state) |>
  summarise(
    mean(pop2011),
    sum(pop2011)
  )

# A tibble: 17 x 3
#> #> #> state      `mean(pop2011)` `sum(pop2011)`
#> #> #> <chr>        <dbl>          <dbl>
#> 1 Baden-Wrttemberg      7580           7580
#> 2 Baden-Württemberg     23680.         7837917
#> 3 Bayern                  23996.        7558677
#> 4 Berlin                 3292365        3292365
#> 5 Brandenburg            18472.        1865632
#> 6 Bremen                  325432.       650863
#> 7 Hamburg                 1706696       1706696
#> 8 Hessen                  22996.        5036121
#> 9 Mecklenburg-Vorpommern  27034.        811005
#> 10 Niedersachsen          24107.        6219515
#> 11 Nordrhein-Westfalen    47465.        18036727
#> 12 Rheinland-Pfalz        25644.        1871995
#> 13 Saarland                NA             NA
#> 14 Sachsen                 27788.        2973351
#> 15 Sachsen-Anhalt          21212.        1993915
#> 16 Schleswig-Holstein      24157.        1739269
#> 17 Th_ringen               29192.        1167692

```

## 9. Collection of exercises

```

df <- df |>
  mutate(state = case_when(
    state == "Baden-Wrttemberg" ~ "Baden-Württemberg",
    state == "Th_ringen" ~ "Thüringen",
    TRUE ~ state
  ))
```

```

df |>
  group_by(state) |>
  summarise(
    mean(pop2011),
    sum(pop2011)
  )
```

# A tibble: 16 x 3

	state	mean(pop2011)	sum(pop2011)
	<chr>	<dbl>	<dbl>
1	Baden-Württemberg	23631.	7845497
2	Bayern	23996.	7558677
3	Berlin	3292365	3292365
4	Brandenburg	18472.	1865632
5	Bremen	325432.	650863
6	Hamburg	1706696	1706696
7	Hessen	22996.	5036121
8	Mecklenburg-Vorpommern	27034.	811005
9	Niedersachsen	24107.	6219515
10	Nordrhein-Westfalen	47465.	18036727
11	Rheinland-Pfalz	25644.	1871995
12	Saarland	NA	NA
13	Sachsen	27788.	2973351
14	Sachsen-Anhalt	21212.	1993915
15	Schleswig-Holstein	24157.	1739269
16	Thüringen	29192.	1167692

```

df |>
  filter(state == "Saarland") |>
  print(n = 100)
```

# A tibble: 47 x 5

	city	status	state	pop2011	rankX
	<chr>	<chr>	<chr>	<dbl>	<dbl>
1	Perl	Commune	Saarland	7775	2003
2	Freisen	Commune	Saarland	8270	1894
3	Großrosseln	Commune	Saarland	8403	1868
4	Nonnweiler	Commune	Saarland	8844	1775
5	Nalbach	Commune	Saarland	9302	1678
6	Wallerfangen	Commune	Saarland	9542	1642
7	Kirkel	Commune	Saarland	10058	1541
8	Merchweiler	Commune	Saarland	10219	1515

## 9. Collection of exercises

9 Nohfelden	Commune	Saarland	10247	1511
10 Friedrichsthal	City	Saarland	10409	1489
11 Marpingen	Commune	Saarland	10590	1461
12 Mandelbachtal	Commune	Saarland	11107	1390
13 Kleinblittersdorf	Commune	Saarland	11396	1354
14 Überherrn	Commune	Saarland	11655	1317
15 Mettlach	Commune	Saarland	12180	1241
16 Tholey	Commune	Saarland	12385	1217
17 Saarwellingen	Commune	Saarland	13348	1104
18 Quierschied	Commune	Saarland	13506	1088
19 Spiesen-Elversberg	Commune	Saarland	13509	1086
20 Rehlingen-Siersburg	Commune	Saarland	14526	996
21 Riegelsberg	Commune	Saarland	14763	982
22 Ottweiler	City	Saarland	14934	969
23 Beckingen	Commune	Saarland	15355	931
24 Losheim am See	Commune	Saarland	15906	887
25 Schiffweiler	Commune	Saarland	15993	882
26 Wadern	City	Saarland	16181	874
27 Schmelz	Commune	Saarland	16435	857
28 Sulzbach/Saar	City	Saarland	16591	849
29 Illingen	Commune	Saarland	16978	827
30 Schwalbach	Commune	Saarland	17320	812
31 Eppelborn	Commune	Saarland	17726	793
32 Wadgassen	Commune	Saarland	17885	785
33 Bexbach	City	Saarland	18038	777
34 Heusweiler	Commune	Saarland	18201	762
35 Püttlingen	City	Saarland	19134	718
36 Lebach	City	Saarland	19484	701
37 Dillingen/Saar	City	Saarland	20253	654
38 Blieskastel	City	Saarland	21255	601
39 St. Wendel	City	Saarland	26220	460
40 Merzig	City	Saarland	29727	392
41 Saarlouis	City	Saarland	34479	323
42 St. Ingbert	City	Saarland	36645	299
43 Völklingen	City	Saarland	38809	279
44 Homburg	City	Saarland	41502	247
45 Neunkirchen	City	Saarland	46172	206
46 Saarbrücken	City	Saarland	175853	43
47 Perl	Commune	Saarland	NA	NA

```
df <- df |>
  filter(!(city == "Perl" & is.na(pop2011)))

df |>
  filter(state == "Saarland") |>
  print(n = 100)
```

```
# A tibble: 46 x 5
  city                      status state   pop2011 rankX
  <chr>                     <chr>  <chr>    <dbl>   <dbl>
1 Nohfelden                Commune Saarland 10247.  1511.
2 Friedrichsthal            City    Saarland 10409.  1489.
3 Marpingen                 Commune Saarland 10590.  1461.
4 Mandelbachtal             Commune Saarland 11107. 1390.
5 Kleinblittersdorf         Commune Saarland 11396. 1354.
6 Überherrn                  Commune Saarland 11655. 1317.
7 Mettlach                   Commune Saarland 12180. 1241.
8 Tholey                     Commune Saarland 12385. 1217.
9 Saarwellingen              Commune Saarland 13348. 1104.
10 Quierschied               Commune Saarland 13506. 1088.
11 Spiesen-Elversberg        Commune Saarland 13509. 1086.
12 Rehlingen-Siersburg       Commune Saarland 14526.  996.
13 Riegelsberg               Commune Saarland 14763. 982.
14 Ottweiler                  City    Saarland 14934. 969.
15 Beckingen                 Commune Saarland 15355. 931.
16 Losheim am See             Commune Saarland 15906. 887.
17 Schiffweiler               Commune Saarland 15993. 882.
18 Wadern                     City    Saarland 16181. 874.
19 Schmelz                    Commune Saarland 16435. 857.
20 Sulzbach/Saar              City    Saarland 16591. 849.
21 Illingen                   Commune Saarland 16978. 827.
22 Schwalbach                 Commune Saarland 17320. 812.
23 Eppelborn                  Commune Saarland 17726. 793.
24 Wadgassen                  Commune Saarland 17885. 785.
25 Bexbach                    City    Saarland 18038. 777.
26 Heusweiler                 Commune Saarland 18201. 762.
27 Püttlingen                 City    Saarland 19134. 718.
28 Lebach                     City    Saarland 19484. 701.
29 Dillingen/Saar              City    Saarland 20253. 654.
30 Blieskastel                City    Saarland 21255. 601.
31 St. Wendel                 City    Saarland 26220. 460.
32 Merzig                      City    Saarland 29727. 392.
33 Saarlouis                  City    Saarland 34479. 323.
34 St. Ingbert                City    Saarland 36645. 299.
35 Völklingen                 City    Saarland 38809. 279.
36 Homburg                     City    Saarland 41502. 247.
37 Neunkirchen                City    Saarland 46172. 206.
38 Saarbrücken                City    Saarland 175853. 43.
39 Perl                        Commune Saarland NA     NA.
```

9. Collection of exercises

1 Perl	Commune	Saarland	7775	2003
2 Freisen	Commune	Saarland	8270	1894
3 Großrosseln	Commune	Saarland	8403	1868
4 Nonnweiler	Commune	Saarland	8844	1775
5 Nalbach	Commune	Saarland	9302	1678
6 Wallerfangen	Commune	Saarland	9542	1642
7 Kirkel	Commune	Saarland	10058	1541
8 Merchweiler	Commune	Saarland	10219	1515
9 Nohfelden	Commune	Saarland	10247	1511
10 Friedrichsthal	City	Saarland	10409	1489
11 Marpingen	Commune	Saarland	10590	1461
12 Mandelbachtal	Commune	Saarland	11107	1390
13 Kleinblittersdorf	Commune	Saarland	11396	1354
14 Überherrn	Commune	Saarland	11655	1317
15 Mettlach	Commune	Saarland	12180	1241
16 Tholey	Commune	Saarland	12385	1217
17 Saarwellingen	Commune	Saarland	13348	1104
18 Quierschied	Commune	Saarland	13506	1088
19 Spiesen-Elversberg	Commune	Saarland	13509	1086
20 Rehlingen-Siersburg	Commune	Saarland	14526	996
21 Riegelsberg	Commune	Saarland	14763	982
22 Ottweiler	City	Saarland	14934	969
23 Beckingen	Commune	Saarland	15355	931
24 Losheim am See	Commune	Saarland	15906	887
25 Schiffweiler	Commune	Saarland	15993	882
26 Wadern	City	Saarland	16181	874
27 Schmelz	Commune	Saarland	16435	857
28 Sulzbach/Saar	City	Saarland	16591	849
29 Illingen	Commune	Saarland	16978	827
30 Schwalbach	Commune	Saarland	17320	812
31 Eppelborn	Commune	Saarland	17726	793
32 Wadgassen	Commune	Saarland	17885	785
33 Bexbach	City	Saarland	18038	777
34 Heusweiler	Commune	Saarland	18201	762
35 Püttlingen	City	Saarland	19134	718
36 Lebach	City	Saarland	19484	701
37 Dillingen/Saar	City	Saarland	20253	654
38 Blieskastel	City	Saarland	21255	601
39 St. Wendel	City	Saarland	26220	460
40 Merzig	City	Saarland	29727	392
41 Saarlouis	City	Saarland	34479	323
42 St. Ingbert	City	Saarland	36645	299
43 Völklingen	City	Saarland	38809	279
44 Homburg	City	Saarland	41502	247
45 Neunkirchen	City	Saarland	46172	206
46 Saarbrücken	City	Saarland	175853	43

## 9. Collection of exercises

```
df |>
  filter(state == "Saarland") |>
  summarise(
    mean(pop2011),
    sum(pop2011)
  )

# A tibble: 1 x 2
`mean(pop2011)` `sum(pop2011)`
<dbl>          <dbl>
1       20850.     959110

df |>
  group_by(city) |>
  mutate(unique_count = n()) |>
  arrange(city, state) |>
  filter(unique_count > 1) |>
  select(city, status, state, starts_with("pop"), unique_count) |>
  print(n = 100)

# A tibble: 23 x 5
# Groups:   city [11]
  city      status        state    pop2011 unique_count
  <chr>     <chr>        <chr>    <dbl>      <int>
1 Bonn     City with County Rights Nordrhein-Westfalen 305765      3
2 Bonn     City with County Rights Nordrhein-Westfalen 305765      3
3 Bonn     City with County Rights Nordrhein-Westfalen 305765      3
4 Brühl    Commune       Baden-Württemberg 13805       2
5 Brühl    City          Nordrhein-Westfalen 43568       2
6 Erbach   City          Baden-Württemberg 13024       2
7 Erbach   City          Hessen        13245       2
8 Fürth   City with County Rights Bayern      115613      2
9 Fürth   Commune       Hessen        10481       2
10 Lichtenau City          Nordrhein-Westfalen 10473       2
11 Lichtenau Commune      Sachsen       7544        2
12 Münster  Commune      Hessen        14071       2
13 Münster  City with County Rights Nordrhein-Westfalen 289576      2
14 Neunkirchen Commune    Nordrhein-Westfalen 13930       2
15 Neunkirchen City         Saarland      46172       2
16 Neuried   Commune      Baden-Württemberg 9383        2
17 Neuried   Commune      Bayern        8277        2
18 Petersberg Commune     Hessen        14766       2
19 Petersberg Commune     Sachsen-Anhalt 10097        2
20 Senden    City          Bayern        21560       2
21 Senden    Commune      Nordrhein-Westfalen 19976      2
22 Staufenberg City         Hessen        8114        2
23 Staufenberg Commune    Niedersachsen 7983        2
```

## 9. Collection of exercises

```
df |>
  group_by(city, state) |>
  mutate(unique_count = n()) |>
  arrange(city, state) |>
  filter(unique_count > 1) |>
  select(city, status, state, starts_with("pop"), unique_count) |>
  print(n = 100)

# A tibble: 3 x 5
# Groups:   city, state [1]
  city   status           state      pop2011  unique_count
  <chr>  <chr>          <chr>      <dbl>      <int>
1 Bonn  City with County Rights Nordrhein-Westfalen 305765      3
2 Bonn  City with County Rights Nordrhein-Westfalen 305765      3
3 Bonn  City with County Rights Nordrhein-Westfalen 305765      3

df <- df |>
  group_by(city, state) |>
  mutate(n_row = row_number()) |>
  filter(n_row == 1) |>
  select(-n_row)

df |>
  group_by(city, state) |>
  mutate(unique_count = n()) |>
  arrange(city, state) |>
  filter(unique_count > 1) |>
  select(city, status, state, starts_with("pop"), unique_count) |>
  print(n = 100)

# A tibble: 0 x 5
# Groups:   city, state [0]
# i 5 variables: city <chr>, status <chr>, state <chr>, pop2011 <dbl>,
#   unique_count <int>

save(df, file = "city_clean.RData")

df <- df |>
  ungroup() |>
  arrange(desc(pop2011)) |>
  mutate(rank = row_number())

df |>
  select(-rankX, -status, -state) |>
  head()

# A tibble: 6 x 3
  city           pop2011  rank
  <chr>          <dbl>    <dbl>
```

## 9. Collection of exercises

```

<chr>          <dbl> <int>
1 Berlin           3292365     1
2 Hamburg          1706696     2
3 München [Munich] 1348335     3
4 Köln [Cologne]   1005775     4
5 Frankfurt am Main 667925      5
6 Düsseldorf [Dusseldorf] 586291     6

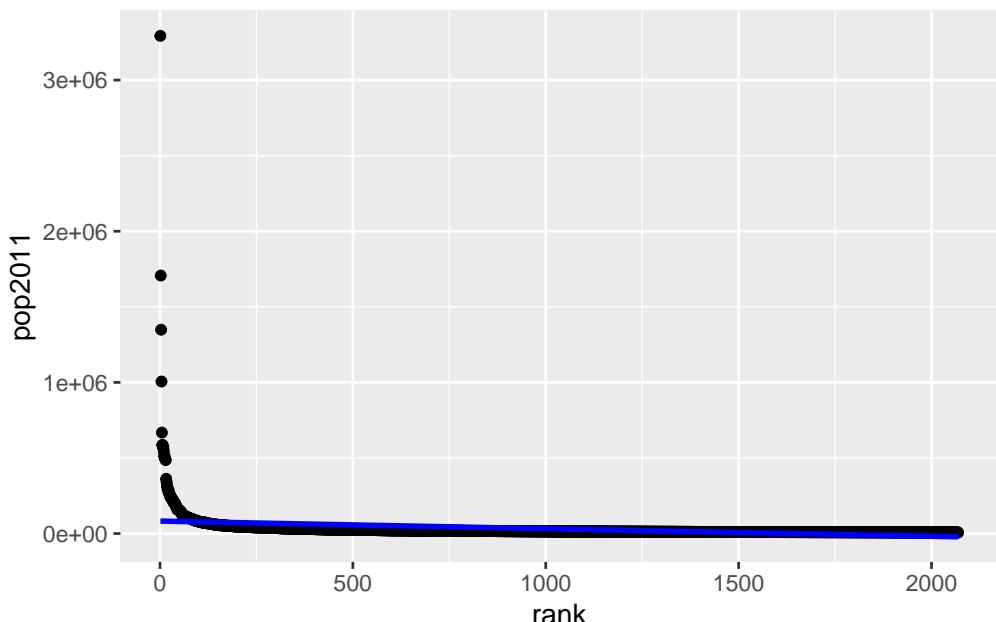
```

```
cor(df$pop2011, df$rank, method = c("pearson"))
```

```
[1] -0.2948903
```

```
ggplot(df, aes(x = rank, y = pop2011)) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE, color = "blue")
```

```
`geom_smooth()` using formula = 'y ~ x'
```



```

df <- df |>
  mutate(lnrank = log(rank)) |>
  mutate(lnpop2011 = log(pop2011))

df |>
  select(city, rank, lnrank, pop2011, lnpop2011) |>
  head()

```

```
# A tibble: 6 x 5
  city              rank  lnrank  pop2011  lnpop2011
  <chr>         <int>  <dbl>    <dbl>      <dbl>
1 Berlin            1     0     3292365    15.0
```

## 9. Collection of exercises

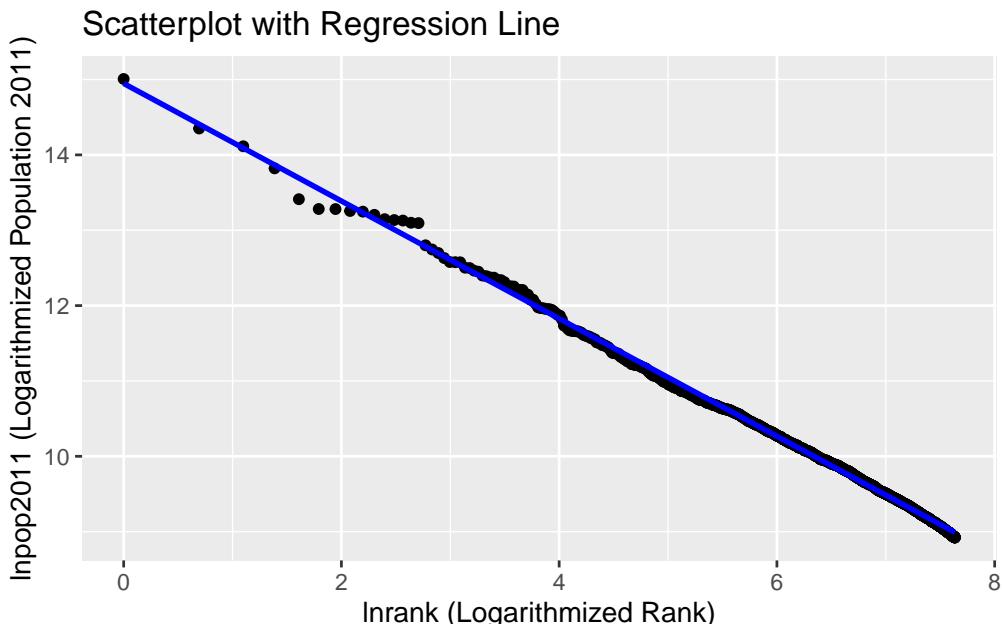
2 Hamburg	2	0.693	1706696	14.4
3 München [Munich]	3	1.10	1348335	14.1
4 Köln [Cologne]	4	1.39	1005775	13.8
5 Frankfurt am Main	5	1.61	667925	13.4
6 Düsseldorf [Dusseldorf]	6	1.79	586291	13.3

```
cor(df$lnpop2011, df$lnrank, method = c("pearson"))
```

[1] -0.9990053

```
ggplot(df, aes(x = lnrank, y = lnpop2011)) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE, color = "blue") +
  labs(
    title = "Scatterplot with Regression Line",
    x = "lnrank (Logarithmized Rank)",
    y = "lnpop2011 (Logarithmized Population 2011)"
  )
```

`geom\_smooth()` using formula = 'y ~ x'



```
zipf <- lm(lnpop2011 ~ lnrank, data = df)
summary(zipf)
```

Call:

lm(formula = lnpop2011 ~ lnrank, data = df)

Residuals:

Min	1Q	Median	3Q	Max
-----	----	--------	----	-----

## 9. Collection of exercises

```
-0.28015 -0.01879  0.01083  0.02005  0.25973

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 14.947859   0.005141    2908 <2e-16 ***
lnrank       -0.780259   0.000766   -1019 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.03454 on 2067 degrees of freedom
Multiple R-squared:  0.998, Adjusted R-squared:  0.998
F-statistic: 1.038e+06 on 1 and 2067 DF,  p-value: < 2.2e-16
```

```
df <- df |>
  mutate(prediction = predict(zipf, newdata = df)) |>
  mutate(pred_pop = exp(prediction))
df |>
  select(city, pop2011, pred_pop) |>
  filter(city == "Regensburg")
```

```
# A tibble: 1 x 3
  city      pop2011 pred_pop
  <chr>     <dbl>    <dbl>
1 Regensburg 135403  134194.
```

```
suppressMessages(pacman::p_unload(tidyverse, haven, janitor, jtools))
```

```
# rmarkdown::render("24-01_dsda.Rmd", "all")
```

```
# knitr::purl(input = "24-01_dsda.Rmd", output = "24-01_dsda_solution.R", documentation = TRUE)
```

**Part IV.**

**Publish**

# 10. Quarto

Verbal and non-verbal communication are crucial in business. This section focuses on writing and publishing texts, excluding aspects like body language and writing skills. I will introduce some tools commonly used by data scientists for writing and publishing their work, such as Markdown, RMarkdown, and Quarto. Unlike applications like Microsoft Word or Apple Pages, these tools use code to generate text. This concept may be unfamiliar to those who grew up after Windows 95; I will provide justification for its use in the following sections.

One notable advantage of a code-based approach to writing text is its seamless integration with version control systems like Git and platforms like GitHub. These tools are essential for most data science collaborations. Mastering them can greatly enhance your efficiency and make your presentations more impactful, even if you are not directly involved in data science.

## What is Quarto?

Quarto, a modern documentation system, is an excellent choice for writing, especially for projects that require rigorous data analysis, visualization, and reproducibility. This tutorial will guide you through producing various forms of text with Quarto. You can write reports, articles, theses, books, websites and many more with Quarto.

## i Quarto and R Markdown

Quarto is a relatively new tool and can be considered as a successor to R Markdown. Most R Markdown documents are compatible with Quarto. However, Quarto offers some improved functionality over R Markdown, which enhances user-friendliness. A detailed overview of the differences and similarities between the two can be found in [this article](#). For an introduction to R Markdown see Section 10.4.

## 10.1. Why Markdown and Quarto?

### 10.1.1. No-code vs. code-based writing applications

Students often use Microsoft Word, Apple Pages, or LibreOffice to write scientific texts. These word processing programs operate on the “What You See Is What You Get” (WYSIWYG) principle, displaying the document layout as you type. While this principle and its corresponding applications are widespread and may seem indispensable to many, this is far from true. Alternatives such as LaTeX, Markdown, R Markdown, and Quarto offer significant advantages. Many professional scientists and publishers prefer these alternatives for good reason. A large number of doctoral theses and scientific papers are authored using LaTeX, and nearly all publishers and editors work with code-based solutions that do not follow the WYSIWYG principle.

With code-based alternatives, layout specifications are either placed at the beginning of the text or embedded within the main text itself. The final document is only visible after converting (also called compiling or rendering) it into a format such as PDF. This may initially seem unusual

and less intuitive than a WYSIWYG interface, but the most intuitive solution is not necessarily the best or simplest. My experience supervising numerous student papers has shown that the intuitive features of MS Word and Pages often become time-consuming over the medium to long term and fail to adequately support users in avoiding errors when writing scientific work. Students who choose code-based applications tend to experience less frustration and greater success—at least, this has been true for the papers I have supervised.

Code-based applications allow writers to focus on the actual writing process, as formatting and adherence to citation rules are largely automated by the software. The necessary initial investment in learning a tool like Quarto quickly pays off, resulting in noticeable improvements in the quality of scientific texts.

In the following subsections, I will first outline typical usage of WYSIWYG applications, then discuss the advantages of code-based text creation using Quarto as an example, and finally explain how to successfully write texts using Quarto.

### 10.1.2. Typical (mis)usage of WYSIWYG applications

The use of traditional word processing software like Microsoft Word or Apple Pages for writing academic papers is pervasive among students. While these programs are user-friendly for everyday writing projects, they create a significant additional workload to meet the demands of academic work.

One of the first problems is integrating literature. Correct formatting according to various citation guidelines is often counter-intuitive, and errors occur easily. This is particularly true if the citation and bibliography functions provided by the software are not used or are used incorrectly. Instead of utilizing external citation managers and investing time to learn how to use them, many students manually create citations and bibliographies, which typically leads to numerous small and sometimes larger errors that could be avoided.

Another weak point in student work is adherence to specific formatting requirements. Academic institutions and journals often require strict adherence to formatting guidelines, including the design of title pages, headers and footers, page margins, and heading hierarchies. Although Word and Pages offer templates and styles, they must be individually adapted for each document and often modified due to minor text changes. Making a format adjustment can become a major effort.

The inclusion of empirical results such as statistical data and graphics presents an additional hurdle. With Word and Pages, the process is frequently manual: research data must be exported from statistical software, saved as images, and then embedded in the document. If the data changes, this time-consuming process must be repeated, significantly increasing the workload and risk of errors.

### 10.1.3. Advantages of Quarto for writing text

Writing academic texts using traditional tools such as MS Word or Pages can be time-consuming and error-prone for students. In the following section, I introduce Quarto (or R Markdown), a modern alternative that offers several advantages:

- **Versatile output formats:** Quarto makes it effortless to generate different output formats. The same text can be rendered as a website (HTML), manuscript (PDF, DOCX), book (EPUB, PDF), or slides (PDF). This flexibility allows you to focus more on the content than the format.

- **Simplified formatting changes:** Specific templates can be used in Quarto, simplifying the process of making formatting changes.
- **Seamless literature integration:** Quarto handles citation rules compliance and integrates seamlessly with citation management systems, enabling researchers to manage literature references and bibliographies more efficiently and consistently than in Word.
- **Easy cross-referencing:** Creating cross-references to sections, tables, and figures is straightforward.
- **Direct data analysis and output generation:** Data analysis and output generation occur directly within Quarto, ensuring that displayed graphics and tables are always up-to-date. This eliminates the need for manual post-processing and guarantees the reproducibility of results.
- **Embedded data visualizations:** Researchers can embed data visualizations directly in the text without manual intermediate steps.
- **Efficient collaboration with version control:** Version control systems like Git make collaboration on academic documents more manageable. Changes can be tracked and integrated without relying on complex and conflict-prone comparison tools.

 Reading recommendation

For those interested, I recommend the online course [Introduction to Reproducible Publications with RStudio](#), which explains explicitly how to work in an empirically reproducible manner. A somewhat more compact introduction is offered by [Bauer and Landesvatter \[2023\]](#), and the authoritative work on the subject is by [Gandrud \[2020\]](#).

## 10.2. Markdown

Markdown is a lightweight markup language with plain-text formatting syntax. It is very popular because it is easy to learn. It's an essential skill for using Quarto effectively. Start by learning enough Markdown to structure your thesis, including headings, lists, links, and code blocks.

You can learn Markdown (not R Markdown!) in 10 minutes. Just go to [www.markdowntutorial.com](#) and work through the interactive lessons. I also recommend the introduction offered in the section [Markdown Basics on quarto.org](#).

## 10.3. Quarto

 Recommended literature

Read [Telford \[2024\]: Enough Markdown to Write a Thesis](#). This resource covers the basics and some advanced Markdown features that are useful for academic writing.

More extensive resources on how to do things with Quarto can be found at [quarto.org](#).

### 10.3.1. Introduction

To set up Quarto on your machine do the following:

- Install R and R Studio.

- Install Quarto as follows:

```
install.packages("quarto")
```

- Install the tinytex package to generate PDF files:

```
install.packages("tinytex")
tinytex::install_tinytex()
```

- It is also advisable to install additional packages that might be needed later:

```
if (!require(pacman)) install.packages("pacman")
pacman::p_load(knitr, rmarkdown, papaja)
```

### Exercise 10.1. First Quarto document

- Open RStudio.
- Select “File” -> “New File” -> “Quarto Document” and then “Create.”
- Save the new file in an empty folder and set this folder as your working directory.
- Click “Render.”
- Visit the Markdown Basics website, add some Markdown to your document, and click “Render” again.
- Click the arrow next to the “Render” button. Here, you can select and generate other file formats. Give it a try.
- Consult the PDF Basics website and supplement your header with the information found there.
- Try citing the paper by Huber and Rust [2016], which you can find here, in your document.
  - Click on “Visual,”
  - Go to the place in the text where you want to cite the paper and select “Insert” -> “Citation.”
  - Search for the paper in the context menu using the corresponding DOI (<https://doi.org/10.1177/1536867X1601600209>) and insert it.
- To quote using APA Version 7 style, write the following in the YAML header:

```
csl: "https://www.zotero.org/styles/apa"
```

- Select a different citation style from www.zotero.org/styles. Then render the document again and observe the differences.

### Exercise 10.2. Quarto cite literature

- a) Create a new Quarto file (*File > New File > Quarto Document*), save the file in an empty folder, and knit it.
- b)
  - Make a new script with *File > New File > R Script*.
  - Go to <https://scholar.google.de/> and search for *osrmtimes*.
  - Click on “cite” and “BibTeX”. Copy and paste everything that you see into your

script and save the script as *lit.bib*. R Studio will ask you if you confirm the file type change. Click yes. Your *lit.bib* file should look like this:

```
@article{huber2016calculate,
  title={Calculate travel time and distance with OpenStreetMap
         data using the Open Source Routing Machine (OSRM)},
  author={Huber, Stephan and Rust, Christoph},
  journal={The Stata Journal},
  volume={16},
  number={2},
  pages={416--423},
  year={2016},
  publisher={SAGE Publications Sage CA: Los Angeles, CA}
}
```

- Add the text “*bibliography: lit.bib*” to your YAML header of your Quarto file so that it looks somehow like that:

```
---
title: "Untitled"
author: "Stephan Huber"
date: "`r Sys.Date()`"
output: html_document
bibliography: lit.bib
---
```

- Now you can cite the OSRMTIME paper with `@huber2016calculate` somewhere in the text of your Quarto file.
- Render the Quarto file and you should see the paper cited and a reference list at the end of the html report.
- You can manipulate the citation style you can specify a CSL (Citation Style Language) file in the YAML header. For example the APA style can be chosen with:

```
csl: "https://www.zotero.org/styles/apa.csl"
```

Many more citation styles can be found on [github.com/citation-style-language](https://github.com/citation-style-language) and on the [Zotero Style Repository](#).

### 10.3.2. Create an APA compliant manuscript using Quarto

To create an APA compliant manuscript, it is recommended to use the *Quarto Extension apaquarto*. The extension that comes with nice templates is hosted on GitHub [here](#). The installation and the usage is described in detail [here](#). Using the template ensures that all APA rules are automatically considered. As APA allows a lot of leeway and every reviewer has specific preferences, `apaquarto` allows for manipulation of a variety of settings. For example, [the language can be changed](#) and the [general style of the document can be modified](#) in the Preamble (YAML header).

Templates for the Fresenius University of Applied Science

If you are a student at Fresenius, feel free to use the templates that I provide. They apply the formating rules of the most study programs offered.

- Quarto template for writing a student paper in English language according the APA 7 Rules
- Quarto template for writing a student paper in German language according the APA 7 Rules
- LaTeX template for writing a thesis at the HS Fresenius in Overleaf or GitHub

## 10.4. R Markdown

Figure 10.1.: Example of an R Markdown file

```

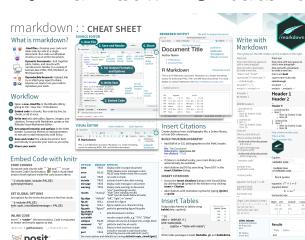
1: <-- 
2: title: "Viridis Demo"
3: output: html_document
4: ---
5: 
6: ````{r include = FALSE}
7: library(viridis)
8: ...
9: 
10: The code below demonstrates two color palettes in the
11: [viridis](https://github.com/simpsoner/viridis) package. Each
12: plot displays a contour map of the Maunga Whau volcano in
13: Auckland, New Zealand.
14: 
15: ## Viridis colors
16: 
17: ````{r}
18: image(volcano, col = viridis(200))
19: 
20: ## Magma colors
21: 
22: ````{r}
23: image(volcano, col = viridis(200, option = "A"))
24: ...
25: 
```

R Markdown provides an authoring framework for data science. You can use a single R Markdown file to transcript your work, run code, and generate high quality reports, books, websites, articles, theses, blogs, and many more (see Figure 10.1).

In contrast to Quarto (see Chapter 10), which is the more recent format, R Markdown is around for some time and hence there are uncountable resources to learn it. For example:

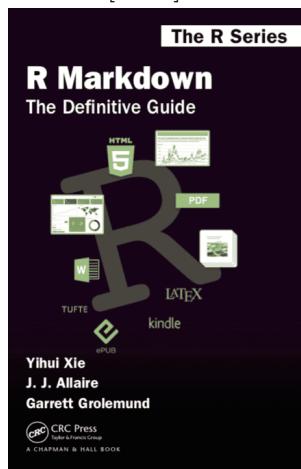
- The [R Markdown Cheatsheet](#) (see Figure 10.2) from Posit offers an overview on the most important features of R Markdown.

Figure 10.2.: R Markdown Cheatsheet from Posit



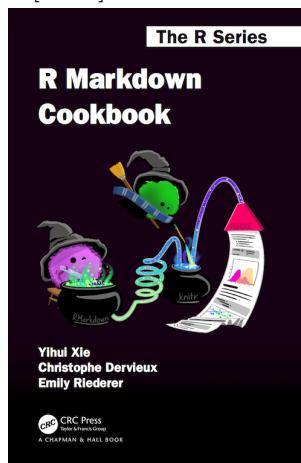
- The book *R Markdown Cookbook* by Xie et al. [2020] (see Figure 10.3) offers an introduction. The [online version of the book](#) is regularly updated and free of costs.

Figure 10.3.: Xie et al. [2020]: R Markdown Cookbook



- The book *R Markdown: The Definitive Guide* by Xie et al. [2018] offers a comprehensive introduction. [The online version of the book](#) is regularly updated and free of costs.

Figure 10.4.: Xie et al. [2018]: R Markdown: The Definitive Guide



### Working directory in R Markdown

The working directory is by default set to the directory that contains the Rmd document. In case you want to use another directory you can do so by changing the working directory with `setwd()`. However, that is not persistent in R Markdown and only works for the current code chunk. After the code chunk has been evaluated, the working directory will be restored to the directory where the Rmd file is placed.

#### Exercise 10.3. Start Markdown and R Markdown

- You can learn Markdown (not R Markdown!) in 10 minutes. Just go to <https://www.markdowntutorial.com> and work through the interactive lessons.
- Now create your first R Markdown file in 3 minutes by doing the following:
  - click in RStudio on *File > New File > R Markdown*
  - click *OK*
  - look for a button entitled *Knit* and click it

- save your file (it will be saved with .Rmd file extension)
- c) Play around with the file. For example, change the output format can you create a word file or a presentation. Play around with the code chunks. Add a picture that you find somewhere online.
- d) Set your working directory to the folder where you have saved your first Rmd-file. Can you come up with a way to generate different output format with just one function.

**Exercise 10.4.** Preparing APA journal articles (`papaja`)

There is an easy way to write a manuscript that follows all the APA rules using the package `papaja` written by two psychologists from Cologne. Please read their [manual](#) and consider their [repository on GitHub](#).

Now, install and load the package:

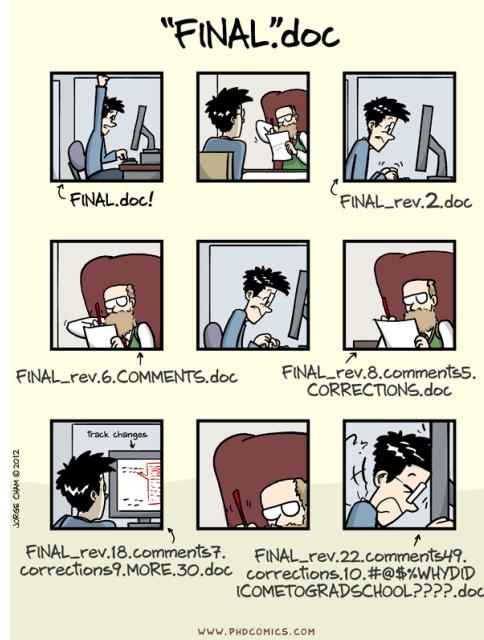
```
install.packages("papaja")
library("papaja")
```

Then, click “File > New File > R Markdown” and choose the “APA-style manuscript” from the section “from template”. Knit the R markdown template and you will have a template for a APA manuscript.

Apart from the obvious adjustments, I recommend to make at least two general adjustments: Change `classoption` to “doc” and `linenumbers` to “no”.

# 11. Collaborating with Git and GitHub

Figure 11.1.: The FINAL.doc problem



Source: [phdcomics.com](http://phdcomics.com)

## 11.1. Introduction

Git is open-source software for version control. It allows developers to track and manage changes to their codebase and files. Users can access a comprehensive history of their project and revert to previous versions of their data if necessary. It helps to overcome the *FINAL.doc* problem depicted in Figure 11.1.



Source: <https://github.com/about as of April 2024>

**GitHub** is an incredibly popular (see statistics in Figure 11.2) online platform that implements Git's capabilities by providing a web interface for collaboration.

While you can use Git and GitHub independently, most developers integrate it with GitHub for enhanced project management and collaboration. This combination helps maintain local and remote copies of a project, facilitating teamwork and data backup as GitHub is sort of a

backup as data loss at your local machine do not matter if you have a remote version saved on GitHub.

Git and GitHub support simultaneous multi-user access, unlike systems that are optimized for single-user like Dropbox.

### 11.2. Install Git

To install the version control system Git, follow the instructions [here](#).

Figure 11.3.: Memorizing six git commands

*The popular approach to version control*



*Source: [DEV Community on GitHub](#)*

Familiarize yourself with Git by using the resources available [here](#). Specifically, work through the resources listed in the box below. Although Git may appear complex, it is generally not too challenging for most users. Many people use Git primarily to track their work and to host and share files conveniently with just a handful of commands. While Git is a robust system with many capabilities, you don't need to memorize all the commands (see Figure 11.3). In fact, you typically use only a few basic ones as shown in Table 11.1.

In the upcoming sections, I will demonstrate some use cases both in the terminal Section 11.3 and within RStudio Section 11.4. In Section 11.5, I show how to contribute to a repository using Git and GitHub.

#### ?

#### Learning resources

Plenty books and tutorial exist that introduce Git and GitHub. I'd like to highlight the following sources:

- The book comprehensive book *Happy Git and GitHub for the useR* by Bryan
- The much shorter book [*Version Control with Git and GitHub*] by Halbritter and Telford
- The online tutorial *How to Use Git/GitHub with R* of David Keyes who explains

in short videos how to setup Git and GitHub in RStudio using among others the `usethis` package.

Table 11.1.: Most important git commands

Git Command	Description
<code>git init</code>	Initialize a new Git repository in the current directory.
<code>git clone &lt;url&gt;</code>	Clone a repository from a remote URL to your local machine.
<code>git add &lt;file&gt;</code>	Add a specific file to the staging area in preparation for committing.
<code>git add .</code>	Add all changed files in the current directory to the staging area.
<code>git commit -m "message"</code>	Commit the staged changes to the repository with a descriptive message.
<code>git status</code>	Display the status of the working directory and staging area.
<code>git push &lt;remote&gt; &lt;branch&gt;</code>	Push committed changes in your local branch to the remote repository.
<code>git pull &lt;remote&gt; &lt;branch&gt;</code>	Pull changes from the remote repository into your current branch and merge them.
<code>git branch &lt;name&gt;</code>	Create a new branch with the specified name.
<code>git checkout &lt;branch&gt;</code>	Switch to another branch and update the working directory.
<code>git merge &lt;branch&gt;</code>	Merge a specified branch into the current branch.

### 11.3. Using Git from the terminal

Figure 11.4.: Three git commands you really need

git add .  
 git commit -m "message"  
 git push

This tutorial will guide you through the basic Git operations using the Bash command line, commonly referred to as the terminal. Essentially, it focuses on the three Git commands illustrated in Figure 11.4.

#### 11.3.1. Configuring Git

Before you start using Git, you need to configure your Git environment. Set your username and email address with these commands:

```
git config --global user.name "Your Name"
git config --global user.email "your.email@example.com"
```

### 11.3.2. Initializing a Repository

To create a new Git repository, use the `git init` command in the directory you want to version control:

```
cd /path/to/a/directory
mkdir my_project
cd my_project
git init
```

In case you are not familiar with using the terminal please consider Table 11.2 where I introduce the most basic commands that we use. For example, with `cd` you can change your directory and with `mkdir` you create a new directory. If you are not familiar with the file system of your computer please read the section [Navigating the file system](#) of Huber [2025]. With `git init` you initialize the directory to be a git repository. This will create a hidden folder “`.git`” in which Git keeps track of all your changes.

#### Most common bash commands

Table 11.2.: Most common bash commands

Bash Command (macOS/Linux)	Windows Command Prompt Equivalent	Description
<code>pwd</code>	<code>cd</code>	Prints the current directory’s path.
<code>ls</code>	<code>dir</code>	Lists all files and directories in the current directory.
<code>cd</code>	<code>cd</code>	Changes the directory.
<code>mkdir</code>	<code>mkdir</code>	Creates a new directory.
<code>rmdir</code>	<code>rmdir</code>	Removes an empty directory.
<code>touch</code>	<code>copy nul</code>	Creates a new empty file or updates an existing file’s timestamp.
<code>rm</code>	<code>del or erase</code>	Removes files. <code>rmdir /s</code> is used for directories.
<code>cp</code>	<code>copy</code>	Copies files or directories.
<code>mv</code>	<code>move</code>	Moves or renames files or directories.
<code>echo</code>	<code>echo</code>	Displays a line of text/string.
<code>cat</code>	<code>type</code>	Concatenates and displays the content of files.
<code>grep</code>	<code>find or findstr</code>	Searches for patterns in files.

### 11.3.3. Staging Changes

To track changes in your repository, you need to stage them using the `git add` command. To stage a single file:

```
git add filename.txt
```

To stage all changes in the directory:

```
git add .
```

#### 11.3.4. Committing Changes

After staging, you can commit it to the repository. A commit records changes to the repository and must include a message describing what changed:

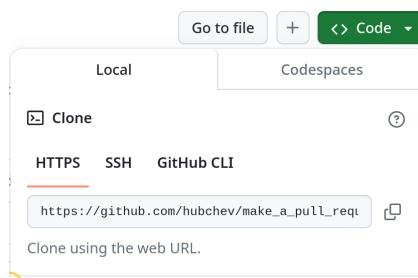
```
git commit -m "A message"
```

#### 11.3.5. Pushing Changes

To share your commits with others or store them in a remote repository (GitHub), use `git push`. A prerequisite here is that you need to be connected to a remote repo. Therefore, you must add a remote repository by copying the URL of the GitHub repo as shown in Figure 11.5. Then you can add the remote repository and push it to the repo with these lines of code:

```
git remote add origin https://github.com/username/repository.git  
git push -u origin main
```

Figure 11.5.: Copy the https URL of your repo



#### 11.3.6. Undo changes

With `git reset` and `git revert` you can go back in time and undo specific changes, respectively. For example, with

```
git log  
git reset --hard <commit_id_hash>
```

you can view the commit history and find the hash identifier of the commit to move the HEAD pointer to that commit. This effectively removes all commits after commit you choose from the current branch's history. Be cautious when using `git reset --hard` as it discards all changes made after the specified commit. Make sure you have backups or are certain you want to discard these changes before proceeding.

With

```
git revert <commit_id_hash>
```

you revert the changes introduced by that commit only. It will create a new commit that undoes the changes made in commit chosen while keeping the other commits that may have followed the chosen commit intact. It's a safer approach compared to `git reset --hard`, as it preserves the commit history and allows you to selectively undo changes without affecting the rest of the commits.

## 11.4. Using Git from RStudio

Integrating Git with RStudio enhances your project management by utilizing version control directly within the IDE. Here's how you can set up and use Git in RStudio using R code.

### 11.4.1. Set up Git in RStudio

First, ensure the `usethis` package is installed and loaded:

```
if (!require(pacman)) install.packages("pacman")
pacman::p_load(usethis)
```

Configure your Git settings in RStudio:

```
use_git_config(user.name = "Your Name", user.email = "Your@email.com")
```

You can change the configuration of your user name and email using the `edit_git_config()` function.

Start a new project in RStudio, which will also initialize a Git repository:

```
create_project("~/Music/")
use_git()
```

After restarting RStudio, you will notice a Git tab in the top right panel, indicating that Git is now active for your project.

### 11.4.2. Connecting RStudio Projects with GitHub repositories

To connect your RStudio project with GitHub, you need a Personal Access Token (PAT) on GitHub. A personal access token (PAT) is required to authenticate with GitHub from Quarto and RStudio. This token allows you to push changes to your repository securely. Follow the instructions to [create a personal access token on GitHub](#). If you haven't one already, you can use the `create_github_token()` function from `usethis` package, and store the PAT securely with `gitcreds_set` from the `gitcreds` package:

```
if (!require(pacman)) install.packages("pacman")
pacman::p_load(usethis gitcreds)
create_github_token()
gitcreds::gitcreds_set()
```

## 11. Collaborating with Git and GitHub

Make sure to note down your token and keep it secure. You'll use this token in RStudio and Quarto to authenticate your GitHub operations.

Now, the procedure depends on whether the project has been initialized on your local machine and you want to create a repo on GitHub, or the repo already exists on GitHub and you want to connect that remote repo with your local PC. Both ways are described below.

### 11.4.2.1. Project exists on RStudio first

After initializing Git in your project, use the `use_github()` function from `usethis` to create a new GitHub repository and connect it directly:

```
use_github()
```

This creates a repo on your GitHub account.

### 11.4.2.2. Project exists on GitHub first

Alternatively, suppose you have created a repository on GitHub first, then start a new project in RStudio using the version control option, specifying your new repository's URL. Just click `File > New Project > Version Control` and then link the GitHub repo by putting the URL into the respective box of the menu. See Figure 11.5 how to get the URL of a repo.

## 11.5. Make a contribution using Git and GitHub

This is a guide for beginners on how to make a contribution using Git and GitHub. If you are looking to make your first contribution, follow the steps below.

### Watch

this [video](#) where I do all the following steps in real time. It takes about 15 minutes.

### 1. Create an account on GitHub

It is free and should just take some minutes.

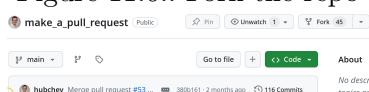
### 2. Install Git

[Here](#) is a tutorial on how to set up Git.

### 3. Fork this repository

Click on the fork button (see Figure 11.6) on the top of this page: [https://github.com/hubchev/make\\_a\\_pull\\_request](https://github.com/hubchev/make_a_pull_request). This will create a copy of this repository in your account.

Figure 11.6.: Fork the repo



### 4. Clone the forked repository

## 11. Collaborating with Git and GitHub

Go to your GitHub account, open the forked repository, click on the code button and then click the *copy to clipboard* icon, see Figure 11.5.

Then, open a terminal and run the following git command:

```
git clone "url you just copied"
```

where “url you just copied” (without the quotation marks) is the url to this repository (your fork of this project). See the previous steps to obtain the url.

For example:

```
git clone https://github.com/this-is-you/make_a_pull_request.git
```

where **this-is-you** is your GitHub username. Here you’re copying the contents of the first-contributions repository on GitHub to your computer.

### 5. Create a branch

Change to the repository directory on your computer (if you are not already there):

```
cd make_a_pull_request
```

Now create a branch using the `git switch` command:

```
git switch -c your-new-branch-name
```

For example:

```
git switch -c add-Stephan-Huber
```

### 6. Make changes.

Now open the `I_am_a_data_scientist.md` file in a text editor. (You find this file in the repository.) Add your name, your GitHub account and the project you are working on. You can put it anywhere in between. Now, save the file.

If you go to the project directory and execute the command `git status`, you’ll see there are changes.

**7. Add changes (staging).** Add those changes to the branch you just created using the `git add` command:

```
git add .
```

**8. Commit changes.** Now commit those changes using the `git commit` command:

```
git commit -m "Add your-name to the list"
```

replacing `your-name` with your name.

**9. Use Git Bash.** Open Git Bash and set your email and your nickname on GitHub:

## 11. Collaborating with Git and GitHub

```
git config --global user.name "FIRST_NAME LAST_NAME"  
git config --global user.email "MY_NAME@example.com"
```

### 10. Push changes to GitHub.

Push your changes using the command `git push`:

```
git push -u origin your-new-branch-name
```

replacing `your-new-branch-name` with the name of the branch you created earlier.

*If you get any errors while pushing that refers to authentication failed something, go to [GitHub's tutorial](#) on generating and configuring an SSH key to your account. Alternatively, you can watch this [YouTube tutorial](#).*

### 11. Submit your changes for review on GitHub.

If you go to your repository on GitHub, you'll see a `Compare & pull request` button. Click on that button.

Now submit the pull request.

Soon I'll be merging all your changes into the main branch of this project. You will get a notification email once the changes have been merged.

Congrats! You just completed the standard *fork -> clone -> edit -> pull request* workflow that you'll often encounter as a contributor!

# 12. Create and host a website

## 12.1. Creating a website with Quarto

This tutorial guides you through creating a simple, yet professional-looking website using Quarto.

### Step W1: Install Quarto

Ensure Quarto is installed on your system. If not, download and install it from [Quarto's official website](#).

### Step W2: Create a website

Follow the tutorial that you find [here](#).

### Step W3: Copy the `_site` directory

After you have rendered your website a directory “`_site`” appears in the project folder that contains your website. Copy all files of that directory to a directory where you want to save your website. Let’s say `my_website`.

In the terminal you can do this with

```
mkdir /home/sthu/my_website/  
cp -r /home/sthu/quarto_website/_site/* /home/sthu/my_website/
```

## 12.2. Hosting the website on GitHub

R Studio and Quarto offers you various ways to publish the website. I explain you a way that worked out well for me.

### Step G1: Create a GitHub account

GitHub will host your thesis website and manage version control for your thesis project. If you don’t already have a GitHub account, you’ll need to create one: Sign up at [GitHub](#).

### Step G2: Create a repository

Create a repository. Name the repo with your username followed by `github.io`. You find a tutorial [here](#).

### Step G3: Obtain a personal access token

A personal access token (PAT) is required to authenticate with GitHub from Quarto and RStudio. This token allows you to push changes to your repository securely. Follow the instructions to [create a personal access token on GitHub](#). Alternatively, you can do the following in R:

```
if (!require(pacman)) install.packages("pacman")
pacman::p_load(usethis)
create_github_token()
```

Make sure to note down your token and keep it secure. You'll use this token in RStudio and Quarto to authenticate your GitHub operations.

#### Step G4: Install and Learn Git

See Section [11.2](#).

#### Step G5: Upload the website to GitHub

Use the Terminal of R Studio. Go to the directory with your website that you have copied in Step W3. Then initiate a git repository on the command line, connect it to the repository created in Steph G2 on GitHub and finally push it:

```
cd /home/sthu/my_website/
echo "# test" >> README.md
git init
git add README.md
git commit -m "first commit"
git branch -M main
git remote add origin https://github.com/test-hsf/test.git
git push -u origin main
```

Alternatively, you can clone a repository, make some changes, and then push those changes back to GitHub. Here are the Bash commands to accomplish this:

```
# Clone the repository
git clone https://github.com/your-username/your-repository.git

# Make changes, here adding a new file as an example
echo "Some content for the new file" > newfile.txt

# Add the new file to the repository
git add newfile.txt

# Commit the changes
git commit -m "Add new file"

# Push the changes back to GitHub
git push origin main
```

# References

- Paul C. Bauer and Camille Landesvatter. Writing a reproducible paper with rstudio and quarto. Technical report, 2023. URL <https://doi.org/10.31219/osf.io/ur4xn>.
- Jennifer Bryan. Happy git and github for the user. URL <https://happygitwithr.com/>.
- John M. Chambers. *Extending R*. CRC Press, 2017.
- Thomas H Davenport and DJ Patil. Data scientist: The sexiest job of the 21st century. *Harvard Business Review*, 90(5):70–76, 2012.
- Christopher Gandrud. *Reproducible research with R and R studio*. Chapman and Hall/CRC, 3 edition, 2020.
- Wickham Hadley. Tidy data. *Journal of Statistical Software*, 59(10):1–23, 2014.
- Aud Halbritter and Richard J Telford. Version control with git and github. URL <https://biostats-r.github.io/biostats/github/>.
- Kieran Healy. *Data Visualization: A Practical Introduction*. Princeton University Press, 2018. URL <https://socviz.co/>.
- Ali Hortacsu and Chad Syverson. The ongoing evolution of US retail: A format tug-of-war. *Journal of Economic Perspectives*, 29(4):89–112, 2015.
- Stephan Huber. How to use R for data science, 2025. URL <https://hubchev.github.io/ds/>.
- Stephan Huber and Christoph Rust. Calculate travel time and distance with openstreetmap data using the open source routing machine (osrm). *The Stata Journal*, 16(2):416–423, 2016. doi: 10.1177/1536867X1601600209.
- Rafael A. Irizarry. *Introduction to Data Science: Data Analysis and Prediction Algorithms With R*. CRC Press, 2022. URL <https://rafalab.github.io/dsbook/>.
- Chester Ismay and Albert Y. Kim. *Statistical inference via data science: A ModernDive into R and the tidyverse*. CRC Press, 2022. URL <https://moderndive.com/>.
- Robert Kabacoff. *Modern Data Visualization with R*. Chapman and Hall/CRC, 2024. URL <https://rkabacoff.github.io/datavis/>.
- John D Kelleher and Brendan Tierney. *Data Science*. MIT Press, 2018.
- Oliver Kirchkamp. Using graphs and visualising data. Technical report, 2018. <https://www.kirchkamp.de/oekonometrie/pdf/gra-p.pdf> (retrieved on 2022/05/20).
- John Muschelli and Andrew Jaffe. Introduction to R for public health researchers. Technical report, GitHub, 2022. URL [https://github.com/muschellij2/intro\\_to\\_r](https://github.com/muschellij2/intro_to_r).
- Danielle Navarro. *Learning Statistics With R*. Version 0.6 edition, 2020. URL <https://learningstatisticswithr.com>.

## References

- Hansjörg Neth. *ds4psy: Data Science for Psychologists*. Social Psychology and Decision Sciences, University of Konstanz, Konstanz, Germany, 2023. URL <https://CRAN.R-project.org/package=ds4psy>. R package (version 0.9.0, October 20, 2022); Textbook at <<https://bookdown.org/hneth/ds4psy/>>.
- Perry Stephenson. Data science practice. Accessed January 30, 2023, 2023. URL <https://datasciencepractice.study/>.
- Richard J Telford. Enough markdown to write a thesis, 8 2024. URL <https://biostats-r.github.io/biostats/quarto/>.
- Måns Thulin. *Modern Statistics With R: From Wrangling and Exploring Data to Inference and Predictive Modelling*. Eos Chasma Press, 2021. URL <https://www.modernstatisticswithr.com/>.
- Edward R. Tufte. *The Visual Display of Quantitative Information*. Graphics Press, 2 edition, 2022.
- William N. Venables, David M. Smith, and R Core Team. *An Introduction to R: Notes on R: A Programming Environment for Data Analysis and Graphics*. Version 4.3.2 (2023-10-31) edition, 2022. URL <http://cran.r-project.org/doc/manuals/R-intro.pdf>. <http://cran.r-project.org/doc/manuals/R-intro.pdf> (retrieved on 2022/04/06).
- Hadley Wickham. The tidyverse style guide, 2024. URL <https://style.tidyverse.org/>.
- Hadley Wickham and Garrett Grolemund. R for data science (2e), 2023. URL <https://r4ds.hadley.nz/>.
- Jeffrey M. Wooldridge. *Introductory Econometrics: A Modern Approach*. South-Western, 2nd edition, 2002.
- Yihui Xie, Joseph J. Allaire, and Garrett Grolemund. *R Markdown: The Definitive Guide*. Chapman and Hall/CRC, 2018.
- Yihui Xie, Christophe Dervieux, and Emily Riederer. *R Markdown Cookbook*. Chapman and Hall/CRC, 2020. available at <https://bookdown.org/yihui/rmarkdown-cookbook>.

# A. Navigating the file system

It is essential to know how R interacts with the file system on your computer. Modern operating systems are incredibly user-friendly and try to hide boring and annoying stuff from the customer. In the following, I will try to give a brief introduction on how to navigate around a computer using a DOS or UNIX shell. If you familiar with that, you can skip this part of the notes.

## A.1. The file system

In this section, I describe the basic idea behind file locations and file paths. Regardless of whether you are using Windows, macOS, or Linux, every file on the computer is assigned a human-readable address, and every address has the same basic structure: it describes a path that starts from a root location, through a series of folders (or directories), and finally ends up at the file.

On a Windows computer, the root is the storage device on which the file is stored, and for many home computers, the name of the storage device that stores all your files is C:. After that comes the folders, and on Windows, the folder names are separated by a backslash symbol \. So, the complete path to this book on my Windows computer might be something like this:

```
C:\Users\huber\Rbook\rcourse-book.pdf
```

On Linux, Unix, and macOS systems, the addresses look a little different, but they are more or less identical in spirit. Instead of using the backslash, folders are separated using a forward slash, and unlike Windows, they do not treat the storage device as being the root of the file system. So, the path on a Mac might be something like this:

```
/Users/huber/Rbook/rcourse-book.pdf
```

That is what we mean by the *path* to a file. The next concept to grasp is the idea of a working directory and how to change it. For those of you who have used command-line interfaces previously, this should be obvious already. But if not, here is what I mean. The working directory is just whatever folder I am currently looking at. Suppose that I am currently looking for files in Explorer (if you are using Windows) or using Finder (on a Mac). The folder I currently have open is my user directory (i.e., C:\Users\huber or /Users/huber). That is my current working directory.

## A.2. Working directory

The next concept to grasp is the idea of a *working directory* and how to change it. For those of you who have used command line interfaces previously, this should be obvious already. But if not, here's what I mean. The working directory is just "whatever folder I'm currently looking at". Suppose that I'm currently looking for files in Explorer (if you're using Windows) or using

### A. Navigating the file system

Finder (on a Mac). The folder I currently have open is my user directory (i.e., C:\Users\huber or /Users/huber). That's my current working directory.

The fact that we can imagine that the program is “in” a particular directory means that we can talk about moving *from* our current location *to* a new one. What that means is that we might want to specify a new location in relation to our current location. To do so, we need to introduce two new conventions. Regardless of what operating system you’re using, we use . to refer to the current working directory, and .. to refer to the directory above it. This allows us to specify a path to a new location in relation to our current location, as the following examples illustrate. Let’s assume that I’m using my Windows computer, and my working directory is C:\Users\huber\Rbook. The table below shows several addresses in relation to my current one:

Absolute path	Relative path
C:\Users\huber	..
C:\Users	..\..
C:\Users\huber\Rbook\source	.\source
C:\Users\huber\nerdstuff	..\nerdstuff

It is quite common on computers that have multiple users to define ~ to be the user’s *home directory*. The home directory on a Mac for the ‘huber’ user is /Users/huber/. And so, not surprisingly, it is possible to define other directories in terms of their relationship to the home directory. For example, an alternative way to describe the location of thercourse-book.pdf“ file on a Mac would be

~\Rbook\rcourse-book.pdf

You can find out your home directory with the `path.expand()` function:

```
path.expand("~/")
```

```
[1] "/home/sthu"
```

Thus, on my machine ~ is an abbreviation for the path /home/sthu.

```
getwd()
```

```
[1] "/home/sthu/Dropbox/hsf/courses/dsr"
```

### A.3. Navigating the file system using the R console

When you want to load or save a file in R it’s important to know what the working directory is. You can find out by using the `getwd()` command. For the moment, let’s assume that I’m using Mac OS or Linux, since things are different on Windows, see section Section A.5. Let’s check the current active working directory:

## A. Navigating the file system

```
getwd()
```

```
[1] "/home/sthu/Dropbox/hsf/courses/dsr"
```

The function `setwd()` allows to change the working directory:

```
setwd("/Users/huber/Rbook/data")
setwd("./Rbook/data")
```

The function `list.files()` lists all the files in that directory:

```
list.files()
```

## A.4. R Studio projects

Setting the working directory repeatedly can be a cumbersome task. Fortunately, R Studio projects can automate this process for you. When you open an R Studio project, the working directory is automatically set to the project directory.

Creating a new project in R Studio is simple. Just click on *File > New Project....*. This will create a directory on your computer with a `*.Rproj` file that can be used to open the saved project at a later date. The newly created directory contains your R code, data files, and other project-related files. By working within projects, all of your files and data are organized in one place, making it easier to share your work with others, reproduce your analyses, and keep track of changes over time.

## A.5. Why do the Windows paths use the back-slash?

Let's suppose I'm using a computer with Windows. As before, I can find out what my current working directory is like this:

```
getwd()
[1] "C:/Users/huber/"
```

R is displaying a Windows path using the wrong type of slash, the back-slash. The answer has to do with the fact that R treats the \ character as *special*. If you're deeply wedded to the idea of specifying a path using the Windows style slashes, then what you need to use two back-slashes \\ whenever you mean \. In other words, if you want to specify the working directory on a Windows computer, you need to use one of the following commands:

```
setwd( "C:/Users/huber" )
setwd( "C:\\Users\\huber" )
```

## B. Operators

### B.1. Assignment:

- `<-` (assignment operator)

### B.2. Arithmetic:

- `+` (addition)
- `-` (subtraction)
- `*` (multiplication)
- `/` (division)
- `^` or `**` (exponentiation)
- `%%` (modulo, remainder)
- `%/%` (integer division)

### B.3. Relational:

- `<` (less than)
- `>` (greater than)
- `<=` (less than or equal to)
- `>=` (greater than or equal to)
- `==` (equal to)
- `!=` or `<>` (not equal to)

### B.4. Logical:

- `&` (element-wise AND)
- `|` (element-wise OR)
- `!` (logical NOT)
- `&&` (scalar AND)
- `||` (scalar OR)

### B.5. Others:

- `%*%` (matrix multiplication)
- `%in%` (checks if an element is in a vector)
- `%>%` or `|>` (pipe operator from the magrittr package)
- `[]`: Extract content from vectors, lists, or data frames.
- `[[ ]]` and `$`: Extract a single item from an object.

## C. Popular functions

### C.1. Help

- `?:` Search R documentation for a specific term.
- `??` Search R help files for a word or phrase.
- `RSiteSearch`: Search search.r-project.org
- `help.start`: Access to html manuals and documentations implemented in R
- `browseVignettes`: view a list of all vignettes associated with your installed packages
- `vignette`: View a specified package vignette, that is, supporting material such as introductions.

### C.2. Package management

- `install.packages`: Installs packages from CRAN.
- `pacman::p_load`: Installs and loads specified R packages.
- `library`: (Install and) loads specified R packages.

### C.3. General

- `setwd`: Sets the working directory to the specified path.
- `rm`: Removes objects (variables) from the workspace.
- `sessionInfo`: Information about the R environment.
- `source`: Executes R code from a file.

### C.4. Tools

- `else`: Execute a block of code if the preceding condition is false.
- `else if`: Specify a new condition to test if the first condition is false.
- `if`: Execute a block of code if a specified condition is true.
- `ifelse`: Check a condition for every element of a vector.

### C.5. Data import

- `c`: Combine values into a vector or list.
- `read.csv`: Reads a CSV file into a data frame.
- `read_dta`: Read Stata dataset.
- `load`: Loads an RData file.

## C.6. Inspect data

- `dim`: Returns the dimensions (number of rows and columns) of a data frame.
- `glimpse`: Provide a concise summary.
- `head`: Returns the first elements.
- `print`: Prints the specified object.
- `names`: Returns the variable names in a data frame.
- `n()` or `nrow()`: Counts the number of observations in a data frame or group of observations.
- `ncol`: Returns the number of columns in a data frame.
- `summary`: Summary statistics.
- `table`: Create a table of counts or cross-tabulation.
- `tail`: Returns the last n elements.
- `unique`: Extracts unique elements from a vector.
- `view`: Opens a viewer for data frames.

## C.7. Graphics

- `abline`: Adds lines to a plot.
- `aes`: Aesthetic mapping in ggplot.
- `facet_wrap`: Creates a grid of faceted plots.
- `geom_hline`: Adds horizontal lines to a ggplot.
- `geom_line`: Adds lines to a ggplot.
- `geom_point`: Adds points to a ggplot.
- `geom_smooth`: Adds a smoothed line to a ggplot.
- `geom_text`: Adds text to a ggplot.
- `geom_vline`: Adds vertical lines to a ggplot.
- `ggsave`: Saves a ggplot to a file.
- `labs`: Adds or modifies plot labels.
- `plot`: Creates a scatter plot.
- `scale_y_reverse`: Reverses the y-axis in a ggplot.
- `stat_smooth`: Adds a smoothed line to a ggplot.
- `theme_classic`: Applies a classic theme to a ggplot.
- `theme_minimal`: Applies a minimal theme to a ggplot.

## C.8. Data management

- `arrange`: Reorder the rows of a data frame.
- `clean_names`: Cleans names of an object (usually a data.frame).
- `complete`: Completes a data frame with all combinations of specified columns.
- `data.frame`: Creates a data frame.
- `distinct`: Removes duplicate rows from a data frame.
- `identical`: Check if two objects are identical.
- `is(na)`: Identify and flag a missing or undefined value (NA).
- `is_tibble`: Check if an object is a tibble.
- `rm`: Removes objects (variables) from the workspace.
- `relocate`: Reorders columns in a data frame.
- `round`: Rounds a numeric vector to the nearest integer.

### C. Popular functions

- `rownames`: Get or set the row names of a matrix-like object.
- `tibble`: Creates a tibble, a modern and tidy data frame.

## C.9. dplyr functions

- `arrange`: Reorder the rows of a data frame.
- `complete`: Completes a data frame with all combinations of specified columns.
- `ends_with`: matches to a specified suffix
- `filter`: Pick observations by their values.
- `first`: Returns the first element.
- `group_by`: Group data by one or more variables.
- `last`: Returns the last element.
- `mutate`: Add new variables or modify existing variables in a data frame.
- `nth`: Returns the nth element.
- `n_distinct`: Returns the number of distinct elements.
- `rename`: Rename variables in a data frame.
- `rename_all`: Renames all variables in a data frame.
- `row_number`: Adds a column with row numbers.
- `rowwise`: Perform operations row by row.
- `select`: Pick variables by their names.
- `select_all`: Selects all columns in a data frame.
- `slice_head`: Selects the top N rows from each group.
- `starts_with`: Select variables whose names start with a certain string.
- `summarise`: Reduce data to a single summary value.

## C.10. Data analysis

- `aggregate`: Apply a function to the data by levels of one or more factors.
- `anti_join`: Return rows from the first data frame that do not have a match in the second data frame.
- `cor`: Computes correlation coefficients.
- `cov`: Computes covariance.
- `diff`: Calculates differences between consecutive elements.
- `get_dupes`: Identify duplicate rows in a data frame (from the janitor package).
- `paste0`: Concatenate vectors after converting to character.
- `predict`: Predict method for model fits.
- `prop.table`: Create a table of proportions.

## C.11. Statistical functions

- `cor()`: Computes correlation coefficients.
- `cov()`: Computes the covariance.
- `exp()`: Exponential function.
- `IQR()`: Computes the interquartile range.
- `kurtosis()`: Computes the kurtosis.
- `log()`: Natural logarithm.
- `mad()`: Computes the mean absolute deviation.

### *C. Popular functions*

- `max()`: Returns the maximum value.
- `mean()`: Calculates the mean.
- `median()`: Computes the median.
- `min()`: Returns the minimum value.
- `quantile()`: Computes sample quantiles.
- `sd()`: Calculates the standard deviation.
- `skewness()`: Calculates the skewness.
- `var()`: Calculates the variance.

## D. Helpful shortcuts

Table D.1.: Different OS, different keys

Key in Windows/Linux	Key in Mac
CTRL	Command Key
Alt	Option Key

Table D.2.: Helpful shortcuts

Action	Shortcut Keys	Description
Run code	Ctrl + Enter	Runs the current line and jumps to the next one, or runs the selected part without jumping further.
	Alt + Enter	Allows running code without moving the cursor to the next line if you want to run one line of code multiple times without selecting it.
	Ctrl + Alt + R	Runs the entire script.
	Ctrl + Alt + B/E	Run the script from the Beginning to the current line and from the current line to the End.
Write code	Alt + (-)	Inserts the assignment operator (<-) with spaces surrounding it.
	Ctrl + Shift + M	Inserts the magrittr/pipe operator (%>%) with spaces surrounding it.
	Ctrl + Shift + C	Comments out code by putting a # in front of each line of marked code of a script.
	Ctrl + Shift + R	Creates a foldable comment section in your code.
Navigating in RStudio	Ctrl + 1	Move focus to editor.
	Ctrl + 2	Move focus to console.
	Ctrl+Tab and	to switch between tabs.
	Ctrl+Shift+Tab	
	Ctrl + Shift + N	Open a new R script.
	Ctrl + w	Close a tab.