

Exercise 1

No Fear of Numbers: Introduction to Quantitative Data Analysis in R

Dr. Susanne de Vogel, Data Science Center, University of Bremen

2 December 2025

Research Question: Is there a Life beyond the PhD?

In this workshop, we will work with one coherent example throughout all four parts. We focus on two outcomes of interest:

- bpsy01: overall life satisfaction (10 point scale)
- blcd06: Children (yes/no)

We will examine how these outcomes are related to different aspects of doctoral researchers' lives and backgrounds:

PhD and Work conditions

- adbi01: Status of the doctorate
- adbi15: Discipline
- bcd17: Emotional Support during PhD
- bwdr12: Perceived scientific pressure
- bemp81: Monthly gross income

Attitudes and Well-Being

- bldc12: Satisfaction with work-life balance
- apar15b: Relationship with parents
- bpsy05: Self-Efficacy

Demographics

- adem01: Gender
- adem02: Age in years
- apar10: Highest vocational degree of parents
- adem03: Country of birth
- blcd01: Relationship status

Before we start: Getting to know RStudio

When you work with R, you usually use **RStudio** — a user-friendly interface that makes it easier to write, run, and organize your R code. RStudio typically shows **four main panels** side by side:

Source (top left)

This is where you **write your code** — in scripts or R Markdown file. this is your main workspace.

Everything you write in here is treated by R as code. To structure your code with headlines or to comment your code, write a `#` in the beginning of the line to tell R that the following content is a comment. Hence, R will not read it as code.

To run your code, click the Run button at the top right of the `source` window and tell R to either `run selected lines` or `run all`.

Console (bottom left)

This is where **R executes your code**. You'll see results printed here.

You can also type commands directly into the Console, but they won't be saved. That's why it's best to write your code in a script or R Markdown file instead.

Environment / History (top right)

This shows which objects (datasets, variables, results, etc.) are **currently loaded in memory**:

- Environment: list your data and variables.
- History: keep track of all commands you've recently used.

Files / Plots / Packages / Help / Viewer (bottom right)

This panel has several **useful tabs**:

- Files: browse your working directory.
- Plots: view your generated plots.
- Packages: install or activate R packages.
- Help: access documentation (e.g., type `?mean` in the Console).
- Viewer: display interactive outputs, such as HTML reports.

Exploring datasets

In this exercise, we will use R tidyverse to explore the dataset. The goal is to learn different ways to *examine the dataset and its variables* as a *preparation for analysis*. You will learn how to

- install and load packages
- read SPSS datasets with tidyverse's `haven` package
- use different functions to inspect the dataset's structure, variable types and classes, labels, and missing values.

Exercise 1.1: Installing and loading packages

Okaaaaaay, let's go!

1. Install `tidyverse` and `haven` if not already installed.

As a first step, we need to make sure all packages we will use for our data exploration are installed and up-to-date: `tidyverse`, but also `haven`.

When using the packages for the *first time*, you could simply install it with `install.packages()`, e.g. `install.packages("tidyverse")`.

Alternatively, you can use a *conditional check* like

```
if (!requireNamespace("tidyverse", quietly = TRUE)) install.packages("tidyverse")
```

`requireNamespace("tidyverse", quietly = TRUE)` returns `TRUE` if the package is available. The `!` means “not”, so the condition becomes `TRUE` only when the package is missing. In that case, `install.packages("tidyverse")` is executed.

This checks automatically whether the package is already installed.

2. Use the `library()` function to load `tidyverse` and `haven` so that we can work with them.

Installing packages need to happen only once, but in order to work with them, you need to load them *every time* of the beginning of a new R session.

Hint: Use the same name of the package, but *without* quotation marks.

Note: What happens when you load `tidyverse`?

- You might see a **Warning** like: “*Package was built under R version 4.x.y*”. This is **not an error**. It only means the package was created with a newer R version. In almost all cases, everything will work fine.
- R shows which **core tidyverse packages** are attached automatically: `ggplot2`, `dplyr`, `tidyr`, `readr`, `purrr`, `tibble`, `stringr`, `forcats` etc.
- Under **Conflicts**, R warns you if functions from different packages have the same name (e.g. `filter()` exists in both `dplyr` and `stats`). By default, the tidyverse version is used. If you want the other one, call it explicitly, e.g. `stats::filter()`.

Exercise 1.2: Opening a dataset from an SPSS file

To work with data in R, we then need to import a dataset as a data frame. The function `read_sav()` from the tidyverse `haven` package helps us read **SPSS files**.

1. For this exercise, we will use the file `01_qa_exploring_data.sav` in the `exercises` folder. Define the `exercise` folder as your working directory with the function `setwd("path")` (Note that R uses `/` instead of a backslash).
2. Choose `mydata1` as a name for your data frame and import it with the following structure: `chosen_name <- read_sav("filename")`.
3. On the right side, in the Environment pane, the data frame `mydata1` should now appear under `Data`. What kind of information about the data frame do you find here?

Exercise 1.3: Exploring the structure of your data frame

To check whether the data are in a tidy structure, tidyverse provides a range of functions that allow us to explore the data frame and its variables.

1. Use the function `glimpse(mydata1)` to take a first look at your data frame. What kind of information do you get from the output?
2. Execute the functions `head(mydata1)` and `tail(mydata1)`. What is the difference to `glimpse()`? Which information do you get?
3. Inspect the data frame in the RStudio Data Viewer using the function `View(mydata1)` (be careful: R is case sensitive, so use V with a capital letter). What can you see here?

Note: `View()` is only for looking at the data, **not for editing**. Any changes you make inside the Data Viewer will not be saved to your dataset. If you want to clean or transform the data, you need to use functions like `mutate()`, `filter()`, or `select()` in your code.

Exercise 1.4: Inspecting variables

To understand what kind of variables we are dealing with, we now inspect a few example variables in more detail.

Your task is to find out which type and class of variable we have, which labels are available, whether there are any missings (NAs).

1. Inspect a categorical variable: `adem01` (gender).

To call not a whole data frame, but a *specific variable* in that data frame, write `mydata1$adem01`. You can read this as: “the variable `adem01` from the dataset `mydata1`”:

- `mydata1` = name of the dataset (the table)
- `$` = “take this specific column from the table”
- `adem01` = name of the variable / column

- 1a. Use `class()` or `glimpse()` to inspect the variable `adem01`. What type of variable is it in R?
- 1b. Use `attr()` with the options “label” and “labels” to check whether a variable label and value labels are stored for `adem01`.
- 1c. Create a simple frequency table for `adem01` using `count()`.

Hint: Use the pipe operator `%>%` to first pass the whole data frame into `count()`, e.g. `mydata1 %>% count(adem01)`, because `count()` needs a data frame as input and cannot be applied directly to a single variable like `mydata1$adem01`.

1d. Check whether there are any missing values using the `sum(is.na())` function.

2. Inspecting a numeric variable: `adem02` (age in years)

- 2a. Use `class()` and `glimpse()` to check the type of the variable `adem02`.**
- 2b. Use `attr()` to check the variable and value labels. Anything noticeable?**
- 2c. Use `summary()` and `count()` to get a quick overview of the range of ages.**

3. Inspecting a key outcome variable: `bpsy01` (overall life satisfaction)

`bpsy01` will later be an important outcome variable in our analyses. Before we use it, we want to understand its type, its label and its range of values.

- 3a. Inspect the class, variable label and value labels of `bpsy01`.**
- 3b. Use `summary()` to look at the distribution and `count` to see, how often each value occurs.**

Take Home Checklist: Exploring a new dataset in R

Step	What to do when you open a new dataset in R	Useful functions / questions
1	Check the overall structure	<code>glimpse()</code> , <code>head()</code> , <code>tail()</code> , <code>View()</code>
2	Understand key variables	<code>class()</code> , <code>attr(..., "label")</code> , <code>attr(..., "labels")</code>
3	Look at typical value ranges	<code>summary()</code> , <code>count()</code>
4	Identify missing and special codes	<code>is.na()</code> , unusual values like -996, -998, -989
5	Take notes	Which variables look promising? Which ones clearly need cleaning?
