

# Learning R for Clinical Trial Data

Maya Gans

2020-03-12



# Contents

<b>1</b>	<b>Preface</b>	<b>5</b>
1.1	The book will walk you through: . . . . .	5
<b>2</b>	<b>Getting Set Up with R</b>	<b>7</b>
2.1	Setup instructions . . . . .	7
2.2	Linux . . . . .	8
2.3	For everyone . . . . .	9
<b>3</b>	<b>RStudio IDE</b>	<b>11</b>
3.1	Getting set up . . . . .	13
3.2	Start RStudio. . . . .	14
3.3	Optional Preferences . . . . .	14
3.4	Organizing your working directory . . . . .	16
3.5	Exercise . . . . .	17
<b>4</b>	<b>Interacting with R</b>	<b>19</b>
4.1	Seeking help . . . . .	20
<b>5</b>	<b>Introduction to R</b>	<b>21</b>
5.1	Learning Objectives . . . . .	21
5.2	Creating objects in R . . . . .	21
5.3	Vectors and data types . . . . .	25
5.4	Subsetting vectors . . . . .	27
5.5	Missing data . . . . .	28
<b>6</b>	<b>Starting With Data</b>	<b>31</b>
<b>7</b>	<b>Manipulating Data</b>	<b>39</b>
<b>8</b>	<b>Visualizing Data</b>	<b>41</b>
<b>9</b>	<b>Creating a Report</b>	<b>43</b>



# Chapter 1

## Preface

This book is geared towards people who are bravely taking the step towards learning R but have yet to even download R on their local machines.

### 1.1 The book will walk you through:

- Getting R Set up on your computer
- Understanding the RStudio IDE
- Using the RStudio IDE
- A gentle introduction to R
- Importing and viewing data
- Manipulating data
- Visualizing data
- Use the generated tables and plots to create a report



## Chapter 2

# Getting Set Up with R

### 2.1 Setup instructions

R and RStudio are separate downloads and installations. R is the underlying statistical computing environment, but using R alone is no fun. RStudio is a graphical integrated development environment (IDE) that makes using R much easier and more interactive. You need to install R before you install RStudio. After installing both programs, you will need to install the **tidyverse** and **haven** packages from within RStudio. Follow the instructions below for your operating system, and then follow the instructions to install **tidyverse** and **haven**.

#### 2.1.1 Windows

If you already have R and RStudio installed Open RStudio, and click on “Help” > “Check for updates”. If a new version is available, quit RStudio, and download the latest version for RStudio. To check which version of R you are using, start RStudio and the first thing that appears in the console indicates the version of R you are running. Alternatively, you can type `sessionInfo()`, which will also display which version of R you are running. Go on the CRAN website and check whether a more recent version is available. If so, please download and install it. You can check [here](#) for more information on how to remove old versions from your system if you wish to do so.

If you don’t have R and RStudio installed:

- Download R from the CRAN website.
- Run the .exe file that was just downloaded
- Go to the RStudio download page
- Under Installers select RStudio x.yy.zzz - Windows Vista/7/8/10 (where x, y, and z represent version numbers)
- Double click the file to install it

Once it's installed, open RStudio to make sure it works and you don't get any error messages.

### 2.1.2 macOS

If you already have R and RStudio installed Open RStudio, and click on “Help” > “Check for updates”. If a new version is available, quit RStudio, and download the latest version for RStudio.

To check the version of R you are using, start RStudio and the first thing that appears on the terminal indicates the version of R you are running. Alternatively, you can type `sessionInfo()`, which will also display which version of R you are running. Go on the CRAN website and check whether a more recent version is available. If so, please download and install it.

If you don't have R and RStudio installed:

- Download R from the CRAN website.
- Select the .pkg file for the latest R version
- Double click on the downloaded file to install R
- It is also a good idea to install XQuartz (needed by some packages)
- Go to the RStudio download page
- Under Installers select RStudio x.yy.zzz - Mac OS X 10.6+ (64-bit) (where x, y, and z represent version numbers)
- Double click the file to install RStudio

Once it's installed, open RStudio to make sure it works and you don't get any error messages.

## 2.2 Linux

Follow the instructions for your distribution from CRAN, they provide information to get the most recent version of R for common distributions. For most distributions, you could use your package manager (e.g., for Debian/Ubuntu run `sudo apt-get install r-base`, and for Fedora `sudo yum install R`), but we don't recommend this approach as the versions provided by this are usually out of date. In any case, make sure you have at least R 3.3.1.

- Go to the RStudio download page
- Under Installers select the version that matches your distribution, and install it with your preferred method (e.g., with Debian/Ubuntu `sudo dpkg -i rstudio-x.yy.zzz-amd64.deb` at the terminal).
- Once it's installed, open RStudio to make sure it works and you don't get any error messages.



## 2.3 For everyone

After installing R and RStudio, you need to install the `tidyverse` and `haven` packages.

After starting RStudio, at the console type:

```
install.packages(c("tidyverse", "haven"))
```

You can also do this by going to Tools -> Install Packages and typing the names of the packages separated by a comma.

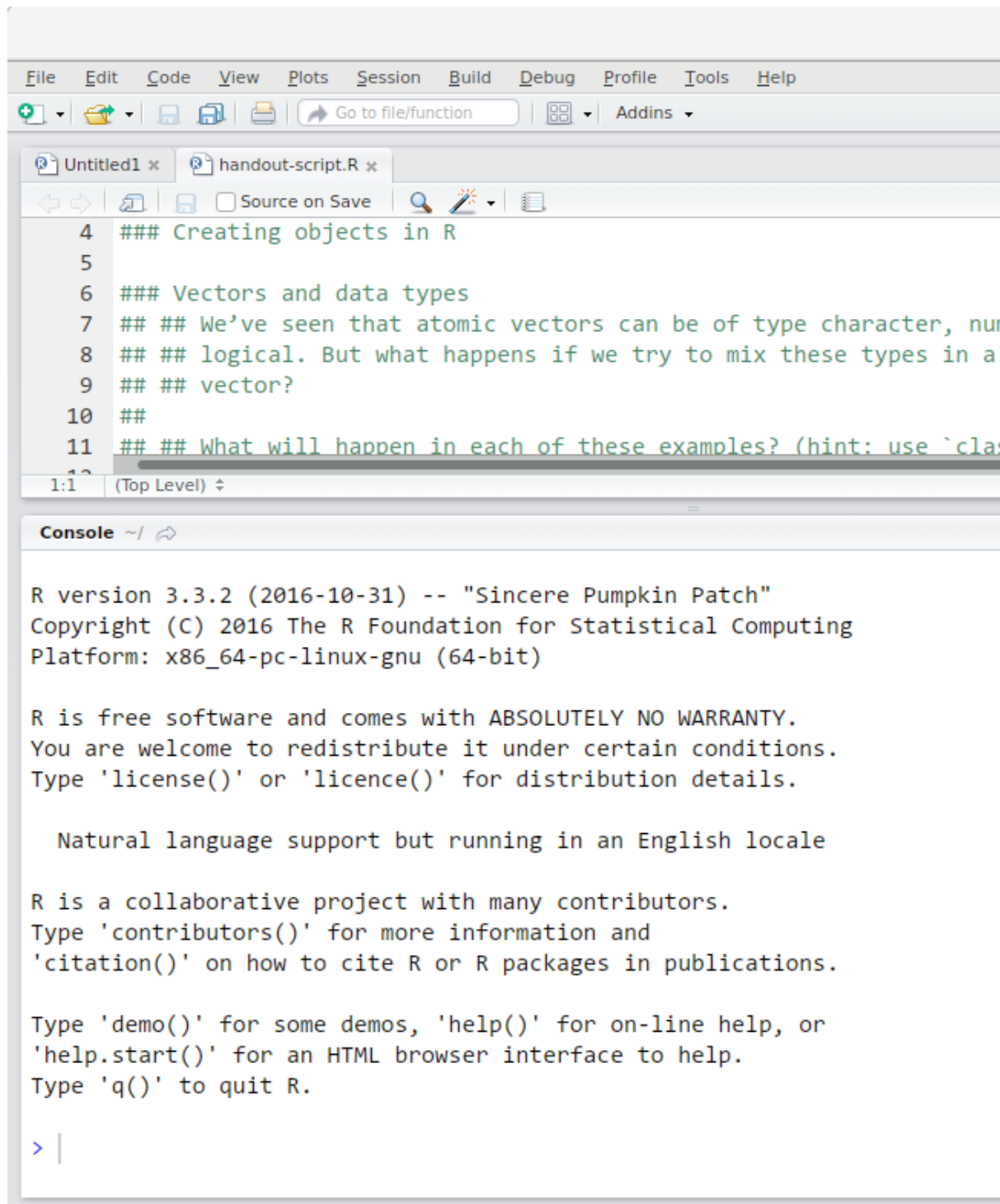


## Chapter 3

# RStudio IDE

The RStudio IDE open-source product is free under the Affero General Public License (AGPL) v3. The RStudio IDE is also available with a commercial license and priority email support from RStudio, Inc.

We will use RStudio IDE to write code, navigate the files on our computer, inspect the variables we are going to create, and visualize the plots we will generate. RStudio can also be used for other things (e.g., version control, developing packages, writing Shiny apps) that will not be covered in this book.



RStudio is divided into 4 “Panels”:

- 1) The Source for your scripts and documents (top-left, in the default layout),
- 2) Your Environment/History (top-right)
- 3) Your Files/Plots/Packages/Help/Viewer (bottom-right)
- 4) The R Console (bottom-left).

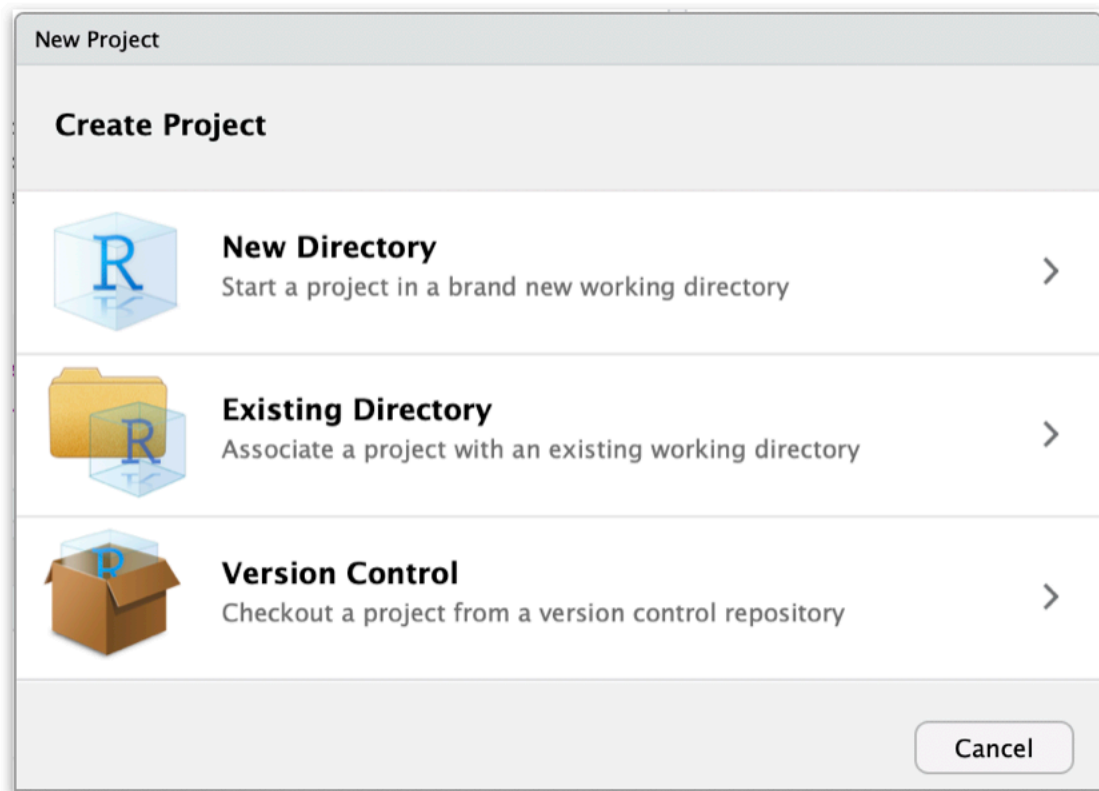
The placement of these panels and their content can be customized (see menu, Tools -> Global Options -> Pane Layout).

One of the advantages of using RStudio is that all the information you need to write code is available in a single window. Additionally, with many shortcuts, autocompletion, and highlighting for the major file types you use while developing in R, RStudio will make typing easier and less error-prone.

## 3.1 Getting set up

It is good practice to keep a set of related data, analyses, and text self-contained in a single folder, called the working directory. All of the scripts within this folder can then use relative paths to files that indicate where inside the project a file is located (as opposed to absolute paths, which point to where a file is on a specific computer). Working this way makes it a lot easier to move your project around on your computer and share it with others without worrying about whether or not the underlying scripts will still work.

RStudio provides a helpful set of tools to do this through its “Projects” interface, which not only creates a working directory for you, but also remembers its location (allowing you to quickly navigate to it) and optionally preserves custom settings and open files to make it easier to resume work after a break. Go through the steps for creating an “R Project” for this tutorial below.



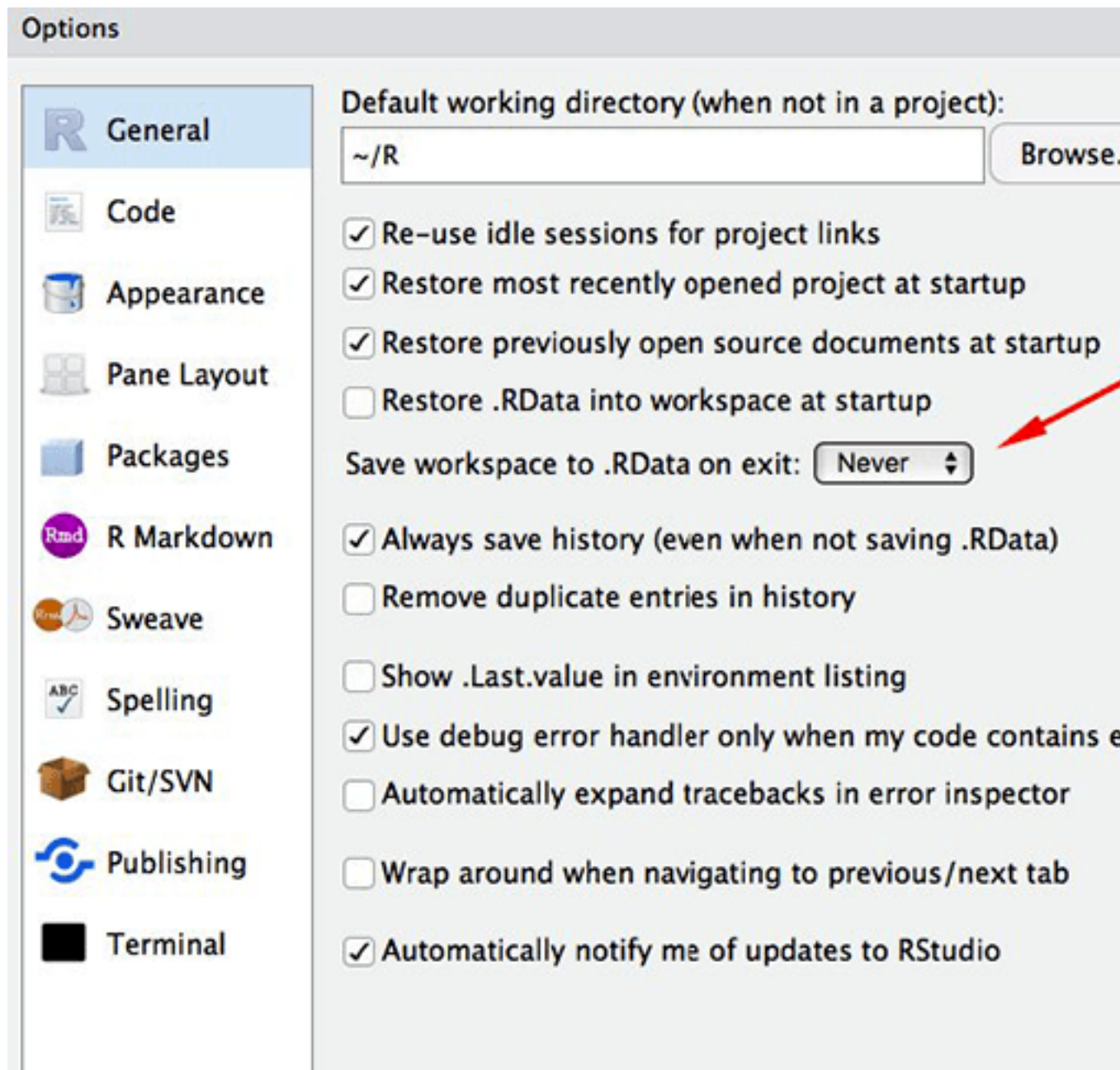
## 3.2 Start RStudio.

- Under the File menu, click on New Project.
- Choose New Directory, then New Project.
- Enter a name for this new folder (or “directory”), and choose a convenient location for it.
- Click on Create Project.

## 3.3 Optional Preferences

RStudio’s default preferences generally work well, but saving a workspace to .RData can be cumbersome, especially if you are working with larger datasets. To turn that off, go to Tools → ‘Global Options’ and select the ‘Never’ option for ‘Save workspace to .RData’ on exit.’ This step is optional, but if you love

something it's sometimes best to let it go.



### 3.4 Organizing your working directory

Using a consistent folder structure across your projects will help keep things organized, and will also make it easy to find/file things in the future. This can be especially helpful when you have multiple projects. In general, you may create directories (folders) for scripts, data, and documents.

- `data_raw/`
- `data/`

Use these folders to store raw data and intermediate datasets you may create for the need of a particular analysis. For the sake of transparency and provenance, you should always keep a copy of your raw data accessible and do as much of your data cleanup and preprocessing programmatically (i.e., with scripts, rather than manually) as possible.

- `report.Rmd`

We will be using an RMarkdown file to create our report. This allows for inline coding with plot and table outputs. We are going to keep the report in the root of our working directory because we are only going to use one file and it will make things easier. Outside of this demonstration you'd most likely create a folder of reports and title them accordingly.

- Additional (sub)directories depending on your project needs (like scripts and functions)

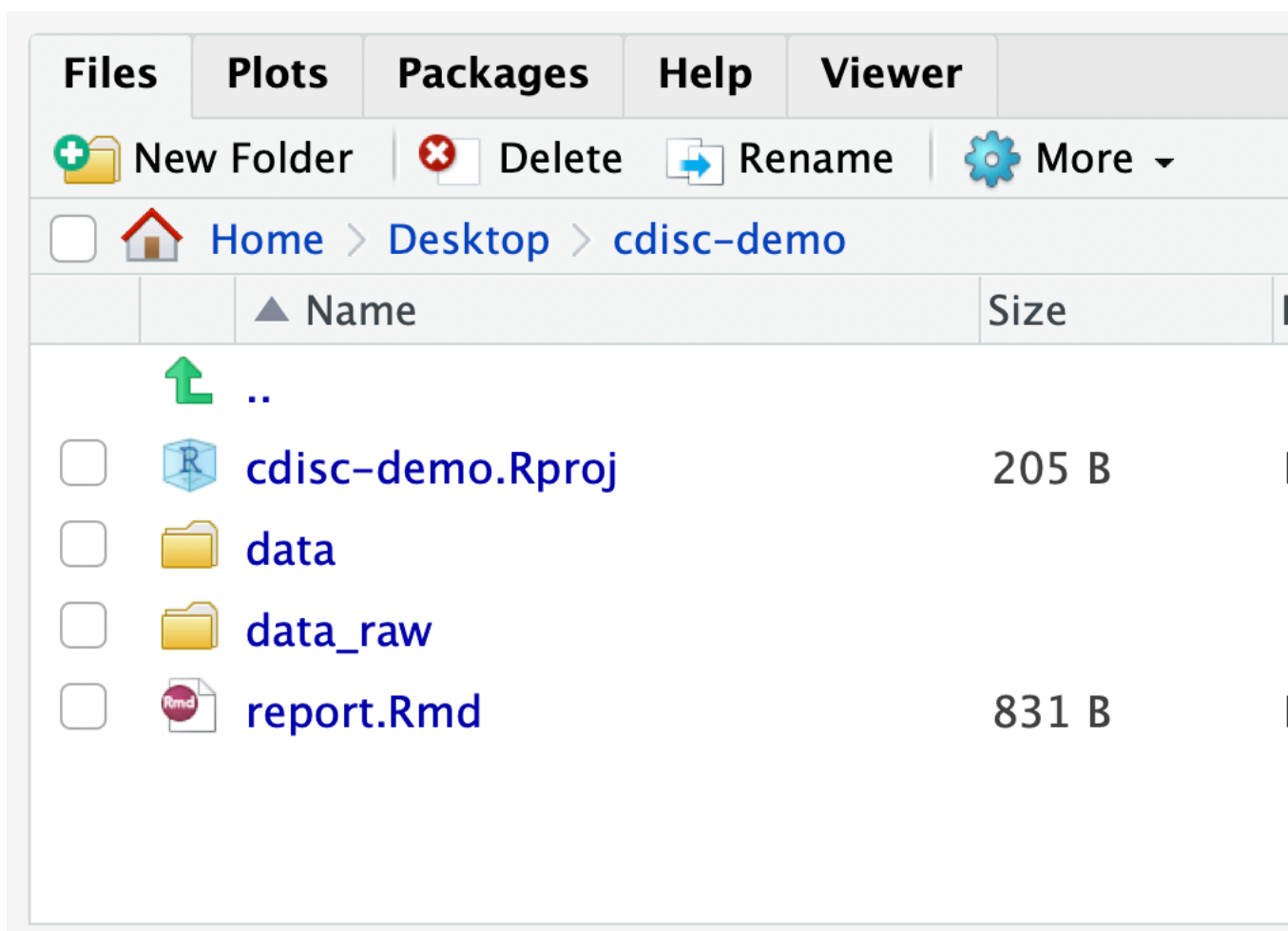
We will need a `data_raw/` folder for our demo project to store our raw `sas7bdat` files, and we will use `data/` for when we learn how to export data as CSV files, and a `report.Rmd` file for our generated report containing figures and tables.

Under the Files tab on the right of the screen, click on New Folder and create a folder named `data_raw` within your newly created working directory (e.g., `~/cdisc-demo/`). (Alternatively, type `dir.create("data_raw")` at your R console.)

Repeat these operations to create a data folder. Under the Files tab you can click New File then RMarkdown and create the `report.Rmd`.

Your working directory should now look like this:





### 3.5 Exercise

Download the ADSL and ADAE from XYZ and add it to your `raw_data` folder



## Chapter 4

# Interacting with R

The basis of programming is that we write down instructions for the computer to follow, and then we tell the computer to follow those instructions. We write, or code, instructions in R because it is a common language that both the computer and we can understand. We call the instructions commands and we tell the computer to follow the instructions by executing (also called running) those commands.

There are two main ways of interacting with R: by using the console or by using script files (plain text files that contain your code). The console pane (in RStudio, the bottom left panel) is the place where commands written in the R language can be typed and executed immediately by the computer. It is also where the results will be shown for commands that have been executed. You can type commands directly into the console and press Enter to execute those commands, but they will be forgotten when you close the session.

Because we want our code and workflow to be reproducible, it is better to type the commands we want in the script editor, and save the script. This way, there is a complete record of what we did, and anyone (including our future selves!) can easily replicate the results on their computer.

RStudio allows you to execute commands directly from the script editor by using the **Ctrl + Enter** shortcut (on Macs, **Cmd + Return** will work, too). The command on the current line in the script (indicated by the cursor) or all of the commands in the currently selected text will be sent to the console and executed when you press Ctrl + Enter. You can find other keyboard shortcuts in this RStudio cheatsheet about the RStudio IDE.

At some point in your analysis you may want to check the content of a variable or the structure of an object, without necessarily keeping a record of it in your script. You can type these commands and execute them directly in the console. RStudio provides the Ctrl + 1 and Ctrl + 2 shortcuts allow you to jump between

the script and the console panes.

If R is ready to accept commands, the R console shows a `>` prompt. If it receives a command (by typing, copy-pasting or sent from the script editor using `Ctrl + Enter`), R will try to execute it, and when ready, will show the results and come back with a new `>` prompt to wait for new commands.

If R is still waiting for you to enter more data because it isn't complete yet, the console will show a `+` prompt. It means that you haven't finished entering a complete command. This is because you have not 'closed' a parenthesis or quotation, i.e. you don't have the same number of left-parentheses as right-parentheses, or the same number of opening and closing quotation marks. When this happens, and you thought you finished typing your command, click inside the console window and press `Esc`; this will cancel the incomplete command and return you to the `>` prompt.

## 4.1 Seeking help

Use the built-in RStudio help interface to search for more information on R functions RStudio help interface. This panel by default can be found at the lower right hand panel of RStudio. As seen in the screenshot, by typing the word "Mean", RStudio tries to also give a number of suggestions that you might be interested in. The description is then shown in the display window.

If you need help with a specific function, let's say `barplot()`, you can type: `?barplot`

If you just need to remind yourself of the names of the arguments, you can use: `args(lm)`

If you are looking for a function to do a particular task, you can use the `help.search()` function, which is called by the double question mark `??`. However, this only looks through the installed packages for help pages with a match to your search request

If you can't find what you are looking for, you can use the [rdocumentation.org](http://rdocumentation.org) website that searches through the help files across all packages available.

Finally, a generic Google or internet search `R <task>` will often either send you to the appropriate package documentation or a helpful forum where someone else has already asked your question.

## Chapter 5

# Introduction to R

### 5.1 Learning Objectives

- Define the following terms as they relate to R: object, assign, call, function, arguments, options.
- Assign values to objects in R.
- Learn how to name objects
- Use comments to inform script.
- Solve simple arithmetic operations in R.
- Call functions and use arguments to change their default options.
- Inspect the content of vectors and manipulate their content.
- Subset and extract values from vectors.
- Analyze vectors with missing data.

### 5.2 Creating objects in R

You can get output from R simply by typing math in the console:

```
3 + 5
```

```
## [1] 8
```

However, to do useful and interesting things, we need to assign values to objects. To create an object, we need to give it a name followed by the assignment operator `<-`, and the value we want to give it:

```
weight_kg <- 55
```

`<-` is the assignment operator. It assigns values on the right to objects on the left. So, after executing `x <- 3`, the value of `x` is 3. The arrow can be read as **3 goes into x**. For historical reasons, you can also use `=` for assignments, but not

in every context. Because of the slight differences in syntax, it is good practice to always use `<-` for assignments.

In RStudio, typing `Alt + -` (push `Alt` at the same time as the `-` key) will write `<-` in a single keystroke in a PC, while typing `Option + -` (push `Option` at the same time as the `-` key) does the same in a Mac.

Objects can be given any name such as `x`, `current_temperature`, or `subject_id`. You want your object names to be explicit and not too long. They **cannot** start with a number (`2x` is not valid, but `x2` is). R is **case sensitive** (e.g., `weight_kg` is different from `Weight_kg`).

There are some names that cannot be used because they are the names of fundamental functions in R (e.g., `if`, `else`, `for`). In general, even if it's allowed, it's best to not use other function names (e.g., `c`, `T`, `mean`, `data`, `df`, `weights`).

If in doubt, check the help to see if the name is already in use. It's also best to avoid dots (`.`) within an object name as in `my.dataset`. There are many functions in R with dots in their names for historical reasons, but because dots have a special meaning in R (for methods) and other programming languages, it's best to avoid them.

It is also recommended to use nouns for object names, and verbs for function names. It's important to be consistent in the styling of your code (where you put spaces, how you name objects, etc.). Using a consistent coding style makes your code clearer to read for your future self and your collaborators. In R, three popular style guides are Google's, Jean Fan's and the tidyverse's. You can install the `lintr` package to automatically check for issues in the styling of your code.

### 5.2.1 Objects vs. variables

What are known as objects in R are known as variables in many other programming languages. Depending on the context, object and variable can have drastically different meanings. However, in this lesson, the two words are used synonymously.

When assigning a value to an object, R does not print anything. You can force R to print the value by using parentheses or by typing the object name:

```
weight_kg <- 55 # doesn't print anything
```

```
(weight_kg <- 55) # but putting parenthesis around the call prints the value of `weight_kg`
```

```
## [1] 55
```

```
weight_kg # and so does typing the name of the object
```

```
## [1] 55
```

Now that R has `weight_kg` in memory, we can do arithmetic with it. For instance, we may want to convert this weight into pounds (weight in pounds is 2.2 times the weight in kg):

```
2.2 * weight_kg
```

```
## [1] 121
```

We can also change an object's value by assigning it a new one:

```
weight_kg <- 57.5
```

```
2.2 * weight_kg
```

```
## [1] 126.5
```

This means that assigning a value to one object does not change the values of other objects. For example, let's store the animal's weight in pounds in a new object, `weight_lb`:

```
weight_lb <- 2.2 * weight_kg
```

and then change `weight_kg` to 100.

```
weight_kg <- 100
```

### 5.2.2 Comments

The comment character in R is `#`, anything to the right of a `#` in a script will be ignored by R. It is useful to leave notes and explanations in your scripts. RStudio makes it easy to comment or uncomment a paragraph: after selecting the lines you want to comment, press at the same time on your keyboard `Ctrl + Shift + C`. If you only want to comment out one line, you can put the cursor at any location of that line (i.e. no need to select the whole line), then press `Ctrl + Shift + C`.

### 5.2.3 Functions and their arguments

Functions are “canned scripts” that automate more complicated sets of commands including operations assignments, etc. Many functions are predefined, or can be made available by importing R packages (more on that later). A function usually takes one or more inputs called arguments. Functions often (but not always) return a value. A typical example would be the function `sqrt()`. The input (the argument) must be a number, and the return value (in fact, the output) is the square root of that number. Executing a function (‘running it’) is called calling the function. An example of a function call is:

```
b <- sqrt(a)
```

Here, the value of `a` is given to the `sqrt()` function, the `sqrt()` function calculates the square root, and returns the value which is then assigned to the object `b`. This function is very simple, because it takes just one argument.

The return ‘value’ of a function need not be numerical (like that of `sqrt()`), and it also does not need to be a single item: it can be a set of things, or even a dataset. We’ll see that when we read data files into R.

Arguments can be anything, not only numbers or filenames, but also other objects. Exactly what each argument means differs per function, and must be looked up in the documentation. Some functions take arguments which may either be specified by the user, or, if left out, take on a default value: these are called options. Options are typically used to alter the way the function operates, such as whether it ignores ‘bad values’, or what symbol to use in a plot. However, if you want something specific, you can specify a value of your choice which will be used instead of the default.

Let’s try a function that can take multiple arguments: `round()`.

```
round(3.14159)
```

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

Here, we’ve called `round()` with just one argument, `3.14159`, and it has returned the value `3`. That’s because the default is to round to the nearest whole number. If we want more digits we can see how to do that by getting information about the round function. We can use `args(round)` to find what arguments it takes, or look at the help for this function using `?round`.

```
args(round)
```

```
## function (x, digits = 0)
## NULL
```

```
?round
```

We see that if we want a different number of digits, we can type `digits = 2` or however many we want.

```
round(3.14159, digits = 2)
```

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

If you provide the arguments in the exact same order as they are defined you don’t have to name them:

```
round(3.14159, 2)
```

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

And if you do name the arguments, you can switch their order:

```
round(digits = 2, x = 3.14159)
```

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



It's good practice to put the non-optional arguments (like the number you're rounding) first in your function call, and to then specify the names of all optional arguments. If you don't, someone reading your code might have to look up the definition of a function with unfamiliar arguments to understand what you're doing.

## 5.3 Vectors and data types

A vector is the most common and basic data type in R, and is pretty much the workhorse of R. A vector is composed by a series of values, which can be either numbers or characters. We can assign a series of values to a vector using the `c()` function. For example we can create a vector of animal weights and assign it to a new object `weight_g`:

```
weight_g <- c(50, 60, 65, 82)
weight_g
```

```
## [1] 50 60 65 82
```

A vector can also contain characters:

```
animals <- c("mouse", "rat", "dog")
animals
```

```
## [1] "mouse" "rat"    "dog"
```

The quotes around “mouse”, “rat”, etc. are essential here. Without the quotes R will assume objects have been created called `mouse`, `rat` and `dog`. As these objects don't exist in R's memory, there will be an error message.

There are many functions that allow you to inspect the content of a vector. `length()` tells you how many elements are in a particular vector:

```
length(weight_g)
```

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

```
length(animals)
```

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

An important feature of a vector, is that all of the elements are the same type of data. The function `class()` indicates the class (the type of element) of an object:

```
class(weight_g)
```

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

```
class(animals)
```

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

The function `str()` provides an overview of the structure of an object and its elements. It is a useful function when working with large and complex objects:

```
str(weight_g)

##  num [1:4] 50 60 65 82
str(animals)

##  chr [1:3] "mouse" "rat" "dog"
```

You can use the `c()` function to add other elements to your vector:

```
weight_g <- c(weight_g, 90) # add to the end of the vector
weight_g

## [1] 50 60 65 82 90

weight_g <- c(30, weight_g) # add to the beginning of the vector
weight_g

## [1] 30 50 60 65 82 90
```

In the first line, we take the original vector `weight_g`, add the value 90 to the end of it, and save the result back into `weight_g`. Then we add the value 30 to the beginning, again saving the result back into `weight_g`.

We can do this over and over again to grow a vector, or assemble a dataset. As we program, this may be useful to add results that we are collecting or calculating.

An **atomic vector** is the simplest R **data type** and is a linear vector of a single type. Above, we saw 2 of the 6 main atomic vector types that R uses: "character" and "numeric" (or "double"). These are the basic building blocks that all R objects are built from. The other 4 atomic vector types are:

- "logical" for TRUE and FALSE (the boolean data type)
- "integer" for integer numbers (e.g., 2L, the L indicates to R that it's an integer)
- "complex" to represent complex numbers with real and imaginary parts (e.g., 1 + 4i) and that's all we're going to say about them
- "raw" for bitstreams that we won't discuss further

You can check the type of your vector using the `typeof()` function and inputting your vector as the argument.

Vectors are one of the many data structures that R uses. Other important ones are lists (`list`), matrices (`matrix`), data frames (`data.frame`), factors (`factor`) and arrays (`array`).

## 5.4 Subsetting vectors

If we want to extract one or several values from a vector, we must provide one or several indices in square brackets. For instance:

```
animals <- c("mouse", "rat", "dog", "cat")
```

```
animals[2]
```

```
## [1] "rat"
```

```
animals[c(3, 2)]
```

```
## [1] "dog" "rat"
```

We can also repeat the indices to create an object with more elements than the original one:

```
more_animals <- animals[c(1, 2, 3, 2, 1, 4)]
more_animals
```

```
## [1] "mouse" "rat" "dog" "rat" "mouse" "cat"
```

R indices start at 1. Programming languages like Fortran, MATLAB, Julia, and R start counting at 1, because that's what human beings typically do. Languages in the C family (including C++, Java, Perl, and Python) count from 0 because that's simpler for computers to do.

### 5.4.1 Conditional subsetting

Another common way of subsetting is by using a logical vector. `TRUE` will select the element with the same index, while `FALSE` will not:

```
weight_g <- c(21, 34, 39, 54, 55)
weight_g[c(TRUE, FALSE, TRUE, TRUE, FALSE)]
```

```
## [1] 21 39 54
```

Typically, these logical vectors are not typed by hand, but are the output of other functions or logical tests. For instance, if you wanted to select only the values above 50:

```
weight_g > 50 # will return logicals with TRUE for the indices that meet the condition
```

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

So we can use this to select only the values above 50

```
weight_g[weight_g > 50]
```

```
## [1] 54 55
```

You can combine multiple tests using `&` (both conditions are true, AND) or `|` (at least one of the conditions is true, OR):

```
weight_g[weight_g < 30 | weight_g > 50]
```

```
## [1] 21 54 55
```

```
weight_g[weight_g >= 30 & weight_g == 21]
```

```
## numeric(0)
```

Here, `<` stands for “less than”, `>` for “greater than”, `>=` for “greater than or equal to”, and `==` for “equal to”. The double equal sign `==` is a test for numerical equality between the left and right hand sides, and should not be confused with the single `=` sign, which performs variable assignment (similar to `<-`).

A common task is to search for certain strings in a vector. One could use the “or” operator `|` to test for equality to multiple values, but this can quickly become tedious. The function `%in%` allows you to test if any of the elements of a search vector are found:

```
animals <- c("mouse", "rat", "dog", "cat")
animals[animals == "cat" | animals == "rat"] # returns both rat and cat
```

```
## [1] "rat" "cat"
```

```
animals %in% c("rat", "cat", "dog", "duck", "goat")
```

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

```
animals[animals %in% c("rat", "cat", "dog", "duck", "goat")]
```

```
## [1] "rat" "dog" "cat"
```

## 5.5 Missing data

As R was designed to analyze datasets, it includes the concept of missing data (which is uncommon in other programming languages). Missing data are represented in vectors as `NA`.

When doing operations on numbers, most functions will return `NA` if the data you are working with include missing values. This feature makes it harder to overlook the cases where you are dealing with missing data. You can add the argument `na.rm = TRUE` to calculate the result while ignoring the missing values.

```
heights <- c(2, 4, 4, NA, 6)
```

```
mean(heights)
```

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

```
max(heights)
```

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

```
mean(heights, na.rm = TRUE)
```

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

```
max(heights, na.rm = TRUE)
```

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

If your data include missing values, you may want to become familiar with the functions `is.na()`, `na.omit()`, and `complete.cases()`. See below for examples.

```
## Extract those elements which are not missing values.
```

```
heights[!is.na(heights)]
```

```
## [1] 2 4 4 6
```

```
## Returns the object with incomplete cases removed. The returned object is an atomic vector of type double.
```

```
na.omit(heights)
```

```
## [1] 2 4 4 6
```

```
## attr("na.action")
```

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

```
## attr("class")
```

```
## [1] "omit"
```

```
## Extract those elements which are complete cases. The returned object is an atomic vector of type double.
```

```
heights[complete.cases(heights)]
```

```
## [1] 2 4 4 6
```

Now that we have learned how to write scripts, and the basics of R's data structures, we are ready to start working with an ADSL dataset!



## Chapter 6

# Starting With Data

Because ADaM datasets are commonly stored as `.sas7bdat` files or `.xpt` files, we will need to install the package `haven`, which enables R to read and write SAS and XPT formats. To install a package we call the function `install.packages`, then to use that library we must call it.

```
# install.packages("haven")
library(haven)
```

Now we can inspect our two datafiles, the ADSL and ADAE taken from XYZ by saving them as objects. We'll use the `head` function to look at the first 6 rows of our data, and the `summary` function to get some basic statistics around each column in each dataframe.

```
ADSL <- read_xpt("data_raw/adsl.xpt")
head(ADSL)
```

```
## # A tibble: 6 x 48
##   STUDYID USUBJID SUBJID SITEID SITEGR1 ARM   TRT01P TRT01PN TRT01A TRT01AN
##   <chr>   <chr>   <chr> <chr>   <chr>   <chr> <chr>   <dbl> <chr>   <dbl>
## 1 CDISCP~ 01-701~ 1015   701     701     Plac~ Plac~       0 Plac~     0
## 2 CDISCP~ 01-701~ 1023   701     701     Plac~ Plac~       0 Plac~     0
## 3 CDISCP~ 01-701~ 1028   701     701     Xano~ Xanom~     81 Xanom~    81
## 4 CDISCP~ 01-701~ 1033   701     701     Xano~ Xanom~     54 Xanom~    54
## 5 CDISCP~ 01-701~ 1034   701     701     Xano~ Xanom~     81 Xanom~    81
## 6 CDISCP~ 01-701~ 1047   701     701     Plac~ Plac~       0 Plac~     0
## # ... with 38 more variables: TRTSDT <date>, TRTEDT <date>, TRTDURD <dbl>,
## #   AVGDD <dbl>, CUMDOSE <dbl>, AGE <dbl>, AGEGR1 <chr>, AGEGR1N <dbl>,
## #   AGEU <chr>, RACE <chr>, RACEN <dbl>, SEX <chr>, ETHNIC <chr>, SAFFL <chr>,
## #   ITTFL <chr>, EFFFL <chr>, COMP8FL <chr>, COMP16FL <chr>, COMP24FL <chr>,
## #   DISCONFL <chr>, DSRAEFL <chr>, DTHFL <chr>, BMIBL <dbl>, BMIBLGR1 <chr>,
## #   HEIGHTBL <dbl>, WEIGHTBL <dbl>, EDUCLVL <dbl>, DISONDT <date>,
```

```
## # DURDIS <dbl>, DURDSGR1 <chr>, VISIT1DT <date>, RFSTDTC <chr>,
## # RFENDTC <chr>, VISNUMEN <dbl>, RFENDT <date>, DCDECOD <chr>, DCSREAS <chr>,
## # MMSETOT <dbl>
```

```
summary(ADSL)
```

```
## STUDYID USUBJID SUBJID SITEID
## Length:254 Length:254 Length:254 Length:254
## Class :character Class :character Class :character Class :character
## Mode :character Mode :character Mode :character Mode :character
##
##
##
## SITEGR1 ARM TRT01P TRT01PN
## Length:254 Length:254 Length:254 Min. : 0.00
## Class :character Class :character Class :character 1st Qu.: 0.00
## Mode :character Mode :character Mode :character Median :54.00
## Mean :44.65
## 3rd Qu.:81.00
## Max. :81.00
##
## TRT01A TRT01AN TRTSDT TRTEDT
## Length:254 Min. : 0.00 Min. :2012-07-09 Min. :2012-08-28
## Class :character 1st Qu.: 0.00 1st Qu.:2013-01-26 1st Qu.:2013-05-12
## Mode :character Median :54.00 Median :2013-06-13 Median :2013-09-27
## Mean :44.65 Mean :2013-06-17 Mean :2013-10-10
## 3rd Qu.:81.00 3rd Qu.:2013-11-07 3rd Qu.:2014-03-15
## Max. :81.00 Max. :2014-09-02 Max. :2015-03-05
##
## TRTDURD AVGDD CUMDOSE AGE
## Min. : 1.00 Min. : 0.00 Min. : 0 Min. :51.00
## 1st Qu.: 46.25 1st Qu.: 0.00 1st Qu.: 0 1st Qu.:70.00
## Median :133.50 Median :54.00 Median : 2322 Median :77.00
## Mean :116.09 Mean :41.54 Mean : 4266 Mean :75.09
## 3rd Qu.:183.00 3rd Qu.:70.10 3rd Qu.: 7810 3rd Qu.:81.00
## Max. :212.00 Max. :78.60 Max. :15417 Max. :89.00
##
## AGEGR1 AGEGR1N AGEU RACE
## Length:254 Min. :1.000 Length:254 Length:254
## Class :character 1st Qu.:2.000 Class :character Class :character
## Mode :character Median :2.000 Mode :character Mode :character
## Mean :2.173
## 3rd Qu.:3.000
## Max. :3.000
##
```



```

##          RACEN          SEX          ETHNIC          SAFFL
##  Min.    :1.00    Length:254    Length:254    Length:254
##  1st Qu.:1.00    Class :character    Class :character    Class :character
##  Median  :1.00    Mode  :character    Mode  :character    Mode  :character
##  Mean    :1.11
##  3rd Qu.:1.00
##  Max.    :6.00
##
##          ITTFL          EFFFLL          COMP8FL          COMP16FL
##  Length:254    Length:254    Length:254    Length:254
##  Class :character    Class :character    Class :character    Class :character
##  Mode  :character    Mode  :character    Mode  :character    Mode  :character
##
##
##
##          COMP24FL          DISCONFL          DSRAEFL          DTHFL
##  Length:254    Length:254    Length:254    Length:254
##  Class :character    Class :character    Class :character    Class :character
##  Mode  :character    Mode  :character    Mode  :character    Mode  :character
##
##
##
##          BMIBL          BMIBLGR1          HEIGHTBL          WEIGHTBL
##  Min.    :13.70    Length:254    Min.    :135.9    Min.    : 34.00
##  1st Qu.:21.90    Class :character    1st Qu.:156.2    1st Qu.: 55.30
##  Median  :24.20    Mode  :character    Median  :162.8    Median  : 66.70
##  Mean    :24.67
##  3rd Qu.:27.30
##  Max.    :40.10
##  NA's    :1
##
##          EDUCLVL          DISONDT          DURDIS          DURDSGR1
##  Min.    : 3.00    Min.    :1998-06-13    Min.    : 2.20    Length:254
##  1st Qu.:12.00    1st Qu.:2008-09-04    1st Qu.: 24.30    Class :character
##  Median  :12.00    Median  :2010-05-05    Median  : 36.25    Mode  :character
##  Mean    :12.75    Mean    :2009-10-09    Mean    : 43.94
##  3rd Qu.:16.00    3rd Qu.:2011-07-16    3rd Qu.: 57.35
##  Max.    :24.00    Max.    :2013-09-16    Max.    :183.10
##
##
##          VISIT1DT          RFSTDTC          RFENDTC          VISNUMEN
##  Min.    :2012-07-06    Length:254    Length:254    Min.    : 4.000
##  1st Qu.:2013-01-17    Class :character    Class :character    1st Qu.: 7.250
##  Median  :2013-06-02    Mode  :character    Mode  :character    Median  :11.000
##  Mean    :2013-06-06
##  3rd Qu.:2013-10-24
##

```

```
## Max.      :2014-08-29                                Max.      :12.000
##
##          RFENDT          DCDECOD          DCSREAS          MMSETOT
## Min.      :2012-09-01    Length:254        Length:254        Min.      :10.00
## 1st Qu.   :2013-05-16    Class :character    Class :character    1st Qu.   :15.00
## Median    :2013-09-30    Mode  :character    Mode  :character    Median   :19.00
## Mean      :2013-10-15                                Mean      :18.14
## 3rd Qu.   :2014-03-20                                3rd Qu.   :22.00
## Max.      :2015-03-05                                Max.      :24.00
##
```

```
ADVS <- read_xpt("data_raw/adae.xpt")
head(ADVS)
```

```
## # A tibble: 6 x 55
##   STUDYID SITEID USUBJID TRTA  TRTAN  AGE AGEGR1 AGEGR1N RACE  RACEN SEX
##   <chr>    <chr>  <chr>   <chr> <dbl> <dbl> <chr>    <dbl> <chr> <dbl> <chr>
## 1 CDISCP~ 701    01-701~ Plac~    0    63 <65      1 WHITE    1 F
## 2 CDISCP~ 701    01-701~ Plac~    0    63 <65      1 WHITE    1 F
## 3 CDISCP~ 701    01-701~ Plac~    0    63 <65      1 WHITE    1 F
## 4 CDISCP~ 701    01-701~ Plac~    0    64 <65      1 WHITE    1 M
## 5 CDISCP~ 701    01-701~ Plac~    0    64 <65      1 WHITE    1 M
## 6 CDISCP~ 701    01-701~ Plac~    0    64 <65      1 WHITE    1 M
## # ... with 44 more variables: SAFFL <chr>, TRTSDT <date>, TRTEDT <date>,
## #   ASTDT <date>, ASTDTF <chr>, ASTDY <dbl>, AENDT <date>, AENDY <dbl>,
## #   ADURN <dbl>, ADURU <chr>, AETERM <chr>, AELLT <chr>, AELLTCD <dbl>,
## #   AEDECOD <chr>, AEPTCD <dbl>, AEHLT <chr>, AEHLTCD <dbl>, AEHLGT <chr>,
## #   AEHLGTCD <dbl>, AEBODSYS <chr>, AESOC <chr>, AESOCCD <dbl>, AESEV <chr>,
## #   AESER <chr>, AESCAN <chr>, AESCONG <chr>, AESDISAB <chr>, AESDTH <chr>,
## #   AESHOSP <chr>, AESLIFE <chr>, AESOD <chr>, AEREL <chr>, AEACN <chr>,
## #   AEOUT <chr>, AESEQ <dbl>, TRTEMFL <chr>, AOCCFL <chr>, AOCCSFL <chr>,
## #   AOCCPFL <chr>, AOCC02FL <chr>, AOCC03FL <chr>, AOCC04FL <chr>,
## #   CQ01NAM <chr>, AOCC01FL <chr>
```

```
summary(ADVS)
```

```
##   STUDYID          SITEID          USUBJID          TRTA
## Length:1191      Length:1191      Length:1191      Length:1191
## Class :character  Class :character  Class :character  Class :character
## Mode  :character  Mode  :character  Mode  :character  Mode  :character
##
##
##
##   TRTAN          AGE          AGEGR1          AGEGR1N
## Min.   : 0.00    Min.   :51.00    Length:1191      Min.   :1.000
## 1st Qu.: 0.00    1st Qu.:69.50    Class :character  1st Qu.:2.000
```

```

## Median :54.00 Median :77.00 Mode :character Median :2.000
## Mean :50.67 Mean :74.82 Mean :2.131
## 3rd Qu.:81.00 3rd Qu.:81.00 3rd Qu.:3.000
## Max. :81.00 Max. :89.00 Max. :3.000
##
## RACE RACEN SEX SAFFL
## Length:1191 Min. :1.000 Length:1191 Length:1191
## Class :character 1st Qu.:1.000 Class :character Class :character
## Mode :character Median :1.000 Mode :character Mode :character
## Mean :1.139
## 3rd Qu.:1.000
## Max. :6.000
##
## TRTSDT TRTEDT ASTDT
## Min. :2012-07-09 Min. :2012-08-28 Min. :1994-04-01
## 1st Qu.:2013-01-25 1st Qu.:2013-05-23 1st Qu.:2013-03-04
## Median :2013-05-20 Median :2013-09-26 Median :2013-06-30
## Mean :2013-06-05 Mean :2013-10-04 Mean :2013-07-09
## 3rd Qu.:2013-10-24 3rd Qu.:2014-02-20 3rd Qu.:2013-11-27
## Max. :2014-09-02 Max. :2015-03-05 Max. :2014-11-03
## NA's :11
## ASTDTF ASTDY AENDT AENDY
## Length:1191 Min. : -6970.00 Min. :2012-08-03 Min. : -2.00
## Class :character 1st Qu.: 15.00 1st Qu.:2013-03-20 1st Qu.: 27.00
## Mode :character Median : 31.50 Median :2013-07-08 Median : 53.00
## Mean : 34.22 Mean :2013-08-03 Mean : 67.14
## 3rd Qu.: 63.00 3rd Qu.:2013-12-14 3rd Qu.:101.25
## Max. : 194.00 Max. :2014-10-31 Max. :211.00
## NA's :11 NA's :473 NA's :473
## ADURN ADURU AETERM AELLT
## Min. : 1.00 Length:1191 Length:1191 Length:1191
## 1st Qu.: 2.00 Class :character Class :character Class :character
## Median : 11.00 Mode :character Mode :character Mode :character
## Mean : 23.84
## 3rd Qu.: 27.00
## Max. :444.00
## NA's :477
## AELLTCD AEDECOD AEPTCD AEHLT
## Min. : NA Length:1191 Min. : NA Length:1191
## 1st Qu.: NA Class :character 1st Qu.: NA Class :character
## Median : NA Mode :character Median : NA Mode :character
## Mean :NaN Mean :NaN
## 3rd Qu.: NA 3rd Qu.: NA
## Max. : NA Max. : NA
## NA's :1191 NA's :1191
## AEHLTCD AEHLGT AEHLGTCD AEBODSYS

```

```

## Min. : NA Length:1191 Min. : NA Length:1191
## 1st Qu.: NA Class :character 1st Qu.: NA Class :character
## Median : NA Mode :character Median : NA Mode :character
## Mean :NaN Mean :NaN
## 3rd Qu.: NA 3rd Qu.: NA
## Max. : NA Max. : NA
## NA's :1191 NA's :1191
## AESOC AESOCCD AESEV AESER
## Length:1191 Min. : NA Length:1191 Length:1191
## Class :character 1st Qu.: NA Class :character Class :character
## Mode :character Median : NA Mode :character Mode :character
## Mean :NaN
## 3rd Qu.: NA
## Max. : NA
## NA's :1191
## AESCAN AESCONG AESDISAB AESDTH
## Length:1191 Length:1191 Length:1191 Length:1191
## Class :character Class :character Class :character Class :character
## Mode :character Mode :character Mode :character Mode :character
##
##
##
## AESHOSP AESLIFE AESOD AEREL
## Length:1191 Length:1191 Length:1191 Length:1191
## Class :character Class :character Class :character Class :character
## Mode :character Mode :character Mode :character Mode :character
##
##
##
## AEACN AEOUT AESEQ TRTEMFL
## Length:1191 Length:1191 Min. : 1.00 Length:1191
## Class :character Class :character 1st Qu.: 2.00 Class :character
## Mode :character Mode :character Median : 4.00 Mode :character
## Mean : 4.53
## 3rd Qu.: 6.00
## Max. :23.00
##
## AOCCFL AOCCSFL AOCCPFL AOCC02FL
## Length:1191 Length:1191 Length:1191 Length:1191
## Class :character Class :character Class :character Class :character
## Mode :character Mode :character Mode :character Mode :character
##
##
##

```

##				
##	AOCC03FL	AOCC04FL	CQ01NAM	AOCC01FL
##	Length:1191	Length:1191	Length:1191	Length:1191
##	Class :character	Class :character	Class :character	Class :character
##	Mode :character	Mode :character	Mode :character	Mode :character
##				
##				
##				



## Chapter 7

# Manipulating Data





## Chapter 8

# Visualizing Data



## Chapter 9

# Creating a Report