

Introduction to R for Life Sciences

João Lourenço, Tania Wyss & Nadine Fournier

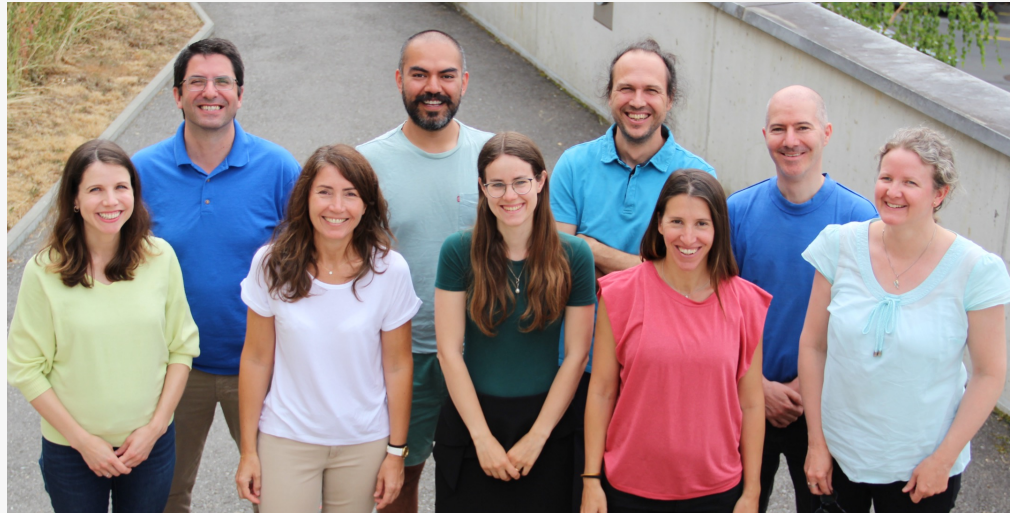
Translational Data Science – Facility

SIB Swiss Institute of Bioinformatics

With slides from Diana Marek, Thomas Junier, Wandrille Duchemin, Leonore Wigger

From: First steps with R in Life Sciences

The Translational Data Science Facility



- Part of the **SIB Swiss Institute of Bioinformatics**
- Located at the AGORA Cancer Research Center in **Lausanne**
- Provides **statistics, bioinformatics and computational expertise** to molecular biology and applied research labs.
- Participates in fundamental and translational research by providing expertise in **data analysis** of single-cell and bulk multi-omics, spatial transcriptomics, flow cytometry, etc

For core facility service inquiry: nadine.fournier@sib.swiss

<https://agora-cancer.ch/scientific-platforms/translational-data-science-facility/>

<https://www.sib.swiss/raphael-gottardo-group>

Tell us about yourself !

Share about yourself and your research,
experience with programming, etc



Photo by National Cancer Institute, Unsplash



Photo by Scott Graham, Unsplash

Course material

1. Website

<https://taniawyss.github.io/intro-to-R/>

The screenshot shows the website interface for 'Analysis of flow cytometry data with R'. The header is blue with the TDS Facility logo, the course title, a search bar, and a GitHub repository link. The main content area is white and divided into three columns. The left column contains a sidebar menu with links to 'Home', 'Intro to R', 'Course schedule', 'Precourse preparations', 'Material', 'Day 1', 'Day 2', and 'Flow cytometry analysis'. The middle column features a 'Home' heading and a detailed description of the course, explaining that it uses R for flexible analysis and pipeline creation, contrasting it with commercial software like FlowJo. It also mentions the course is proposed by the Translational Data Science Facility of the SIB Swiss Institute of Bioinformatics in Lausanne. The right column contains a 'Table of contents' section with links to 'Prerequisite', 'Learning outcomes', 'General learning outcomes', 'Exercises', and 'Asking questions'.

Analysis of flow cytometry data with R

Home

Life scientists often use commercial software such as FlowJo or the OMIQ platform to analyze flow cytometry data. These tools are useful for initial and basic analysis, but do not allow for more advanced or flexible analyses, nor for the establishment of pipelines and reports. On the other hand, R is statistical software that allows for very flexible analysis, customizable pipeline creation and generation of reports.

The “Analysis of flow cytometry data with R” training that is proposed will focus on using R to analyze flow cytometry data. Flow cytometry data that can be analyzed with R includes classical multicolor flow cytometry, spectral flow cytometry, and CyTOF. This course will teach experts in flow cytometry how to run data analysis, develop pipelines and create reports using the open-source R software.

This course is proposed by the [Translational Data Science Facility](#) of the SIB Swiss Institute of Bioinformatics in Lausanne.

Table of contents

- Prerequisite
- Learning outcomes
 - General learning outcomes
- Exercises
- Asking questions

2. Ask us questions!

Outline & Schedule

Morning

01

**Introduction to R and the RStudio environment,
working with scripts files**

Exercises

(9:00 – 10:30)

10:30 – 10:50 Coffee break

02

Syntax, data types and structures, importing data

Exercises

(10:50 -12:00)

12:00 – 13:00 Lunch break

Outline & Schedule

Afternoon

03

Graphics

Exercises

(13:00 – 15:30)

15:30 -15:50 Coffee break

04

04

Statistics

Exercises

(15:50 – 16:50)

16:50 - 17:00 Feedback and end of day

Course Content

R is vast and can't be learned in one day. The scope of this course is to:

- Give you a basic understanding of concepts behind R
- Allow you to import and manipulate data in R
- Show you how to create your first plots

This course is only the first step in your  journey!

01

Introduction to R and the RStudio environment

What is R ?

- R is a **programming language** and an **environment** for statistical computation and graphics.
- A simple **development environment** with a **console** and a **text editor**
- Facilities for **data import, manipulation and storage**
- Functions for **calculations on vectors and matrices**
- Large collections of **data analysis tools**
- **Graphical tools**

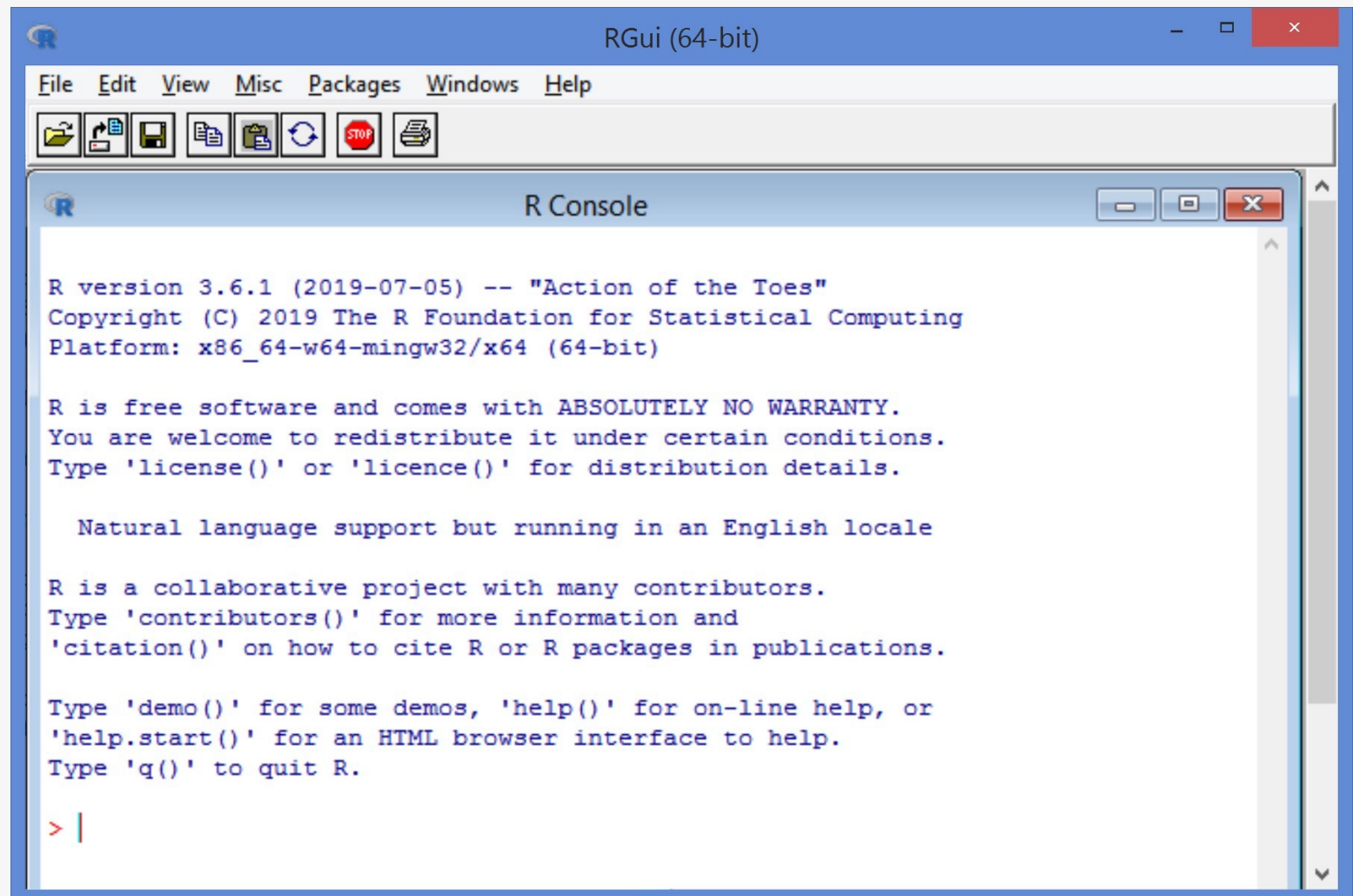
<https://www.r-project.org/>

R's user community

- Group of **core developers** who **maintain** and **upgrade** the basic R installation. New version every 6 months.
- Anyone can contribute with **add-on packages** which provide additional functionality (thousands of such packages available) and **help** for each function.
- **Online help**
 - in user group forums, *eg*:
<https://stat.ethz.ch/mailman/listinfo/r-help>
<http://stackoverflow.com/questions/tagged/r>
 - in countless online tutorials, books, blogs

RGui (R Graphical user interface)

- Together with the programming language, a (minimal) graphical user interface is installed.



R Combined with RStudio



<https://posit.co/products/open-source/rstudio/>

RStudio is an integrated development environment (IDE), designed to help you be more productive with R

It includes:

- A console
- A syntax-highlighting editor that supports direct code execution
- Tools for viewing the workspace and the history
- A file explorer, a package explorer, plot and help display areas

We suggest RStudio as a more powerful, more comfortable alternative to the RGUI

RStudio interface

The screenshot shows the RStudio interface with the following components labeled:

- Editor (scripts)**: The central area for writing R code. It contains a script named `first_script.R` with the following content:

```
1 ##### My first script #####
2 ##### October 2017 #####
3
4
5 # list workspace
6 ls()
7
8 # Reset R's brain
9 rm(list=ls())
10
11 # check wd
12 getwd()
13
14 #set wd
15 setwd("~/Users/dmarek/EducationSIB/Courses_2016/First_Steps_R_June2016/R_intro_course")
16
17 # confirm wd
18 getwd()
19
20 #load packages if needed (to do every time you launch your R session)
21 #library("boot")
22 #library("lattice")
23
24 2:26 10 October 2017 R Script
```
- Workspace (Environment and History)**: The top right pane showing the current environment. It lists objects in the workspace:

Object	Class	Attributes
<code>mice_data</code>	data.frame	50 obs. of 3 variables
<code>mice_weight_HFD</code>	data.frame	29 obs. of 3 variables
<code>mids</code>	num	[1:2, 1] 0.7 1.9
<code>mean_weight_diet</code>	num	[1:2(1d)] 28.7 37.1
<code>mean_weight_genotype</code>	num	[1:2(1d)] 33.7 33.4
<code>n_weight_diet</code>	int	[1:2(1d)] 21 29
<code>n_weight_genotype</code>	int	[1:2(1d)] 24 26
<code>sd_weight_diet</code>	num	[1:2(1d)] 2.61 5
<code>sd_weight_genotype</code>	num	[1:2(1d)] 4.69 6.92
- Console (or terminal)**: The bottom left pane showing the R console output. It displays the R license notice, the R project description, and the workspace loading message:

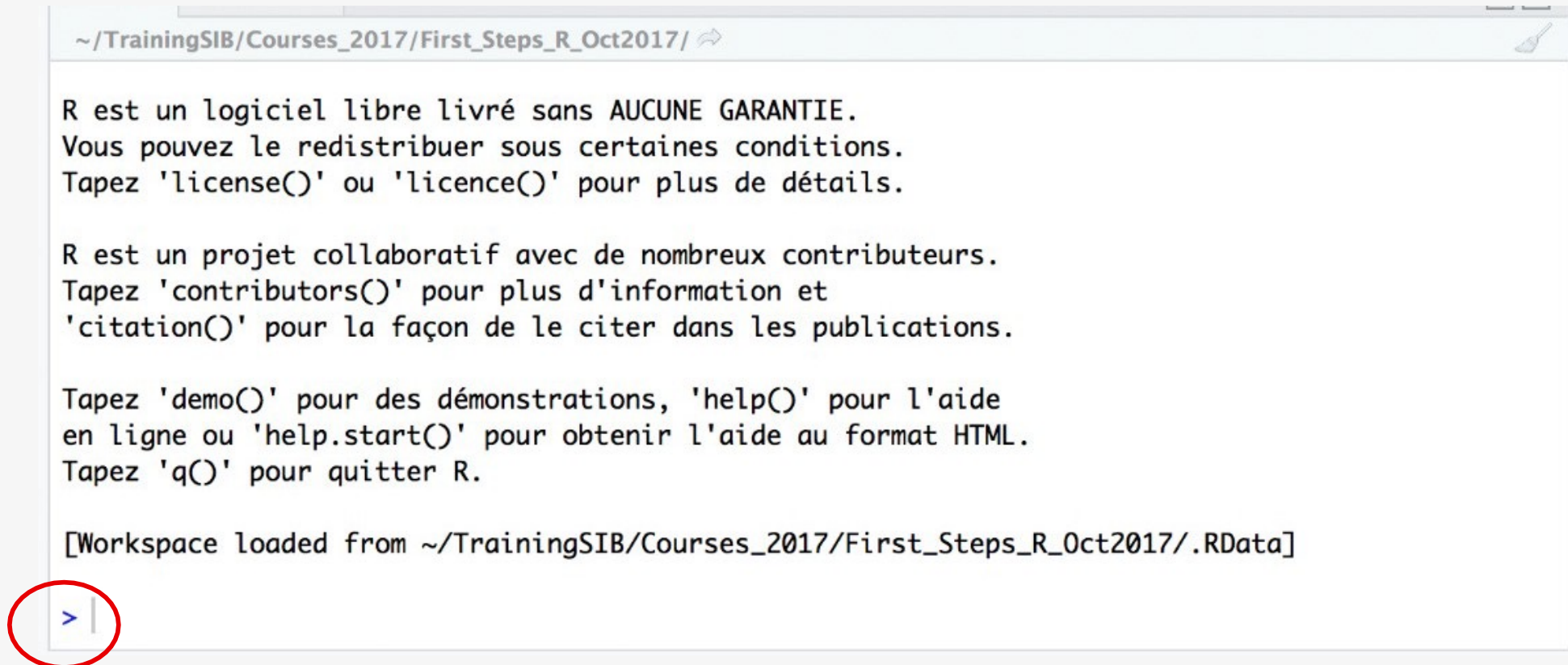
```
R est un logiciel libre livré sans AUCUNE GARANTIE.
Vous pouvez le redistribuer sous certaines conditions.
Tapez 'license()' ou 'licence()' pour plus de détails.

R est un projet collaboratif avec de nombreux contributeurs.
Tapez 'contributors()' pour plus d'information et
'citation()' pour la façon de le citer dans les publications.

Tapez 'demo()' pour des démonstrations, 'help()' pour l'aide
en ligne ou 'help.start()' pour obtenir l'aide au format HTML.
Tapez 'q()' pour quitter R.

[Workspace loaded from ~/TrainingSIB/Courses_2017/First_Steps_R_Oct2017/.RData]
> |
```
- File explorer, plots, packages, help**: The bottom right pane showing the file explorer. It displays the file structure of the current project, including a list of CSV files and their sizes and modification dates.

Console: The Command Line



The screenshot shows an R console window with a title bar indicating the path `~/TrainingSIB/Courses_2017/First_Steps_R_Oct2017/`. The console displays the following text:

```
R est un logiciel libre livré sans AUCUNE GARANTIE.  
Vous pouvez le redistribuer sous certaines conditions.  
Tapez 'license()' ou 'licence()' pour plus de détails.  
  
R est un projet collaboratif avec de nombreux contributeurs.  
Tapez 'contributors()' pour plus d'information et  
'citation()' pour la façon de le citer dans les publications.  
  
Tapez 'demo()' pour des démonstrations, 'help()' pour l'aide  
en ligne ou 'help.start()' pour obtenir l'aide au format HTML.  
Tapez 'q()' pour quitter R.  
  
[Workspace loaded from ~/TrainingSIB/Courses_2017/First_Steps_R_Oct2017/.RData]  
> |
```

The prompt `> |` is circled in red.

The prompt ">" indicates that R is waiting for you to type a command

Try it out...

Type the following at the command prompt:

Simple calculations

```
> 1 + 1
```

Assign values to a variable names

```
> x <- 128.5
```

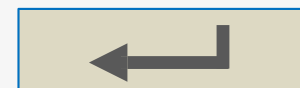
Display content of variables

```
> x
```

Pre-defined functions

```
> abs(-11)
```

**After each command,
hit the return key.**



This causes R to execute it.

Note the assignment operator `<-` with which we can keep values in the memory, by assigning a value and a name to a variable and store it in the session's memory.

We can use either `<-` or `=` to assign values to an object

Stick to one for consistency.

02

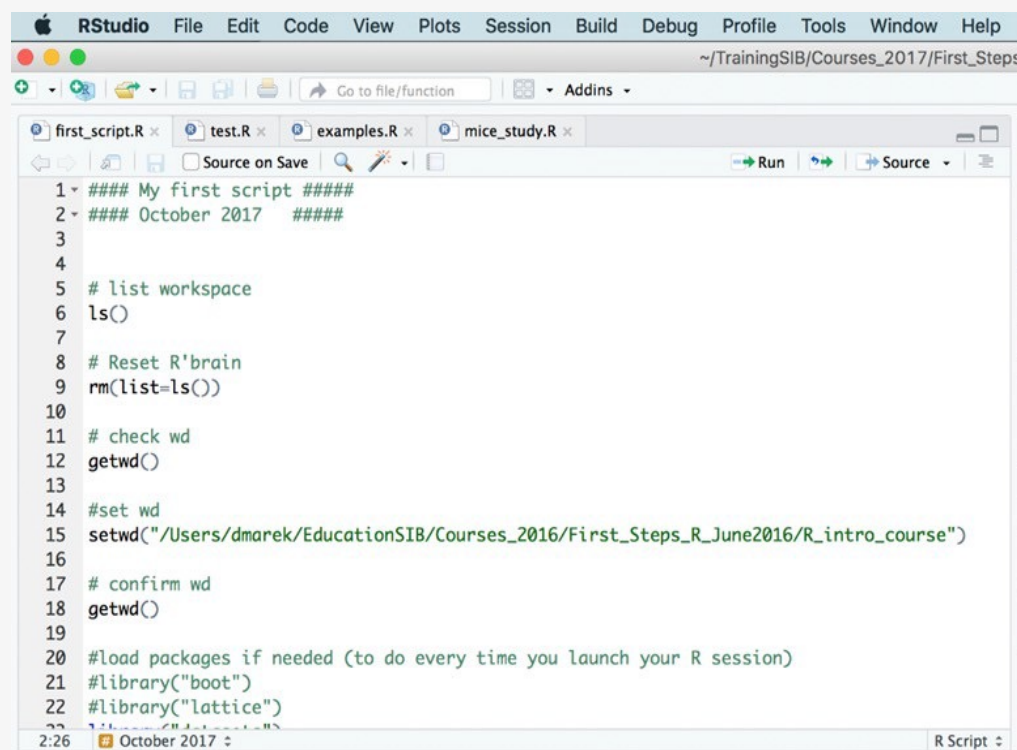
Working with script files

Editor: Write code to a script file

A script is a file that contains commands to be executed in succession.

Write your code into a script and save it

- to have documentation later of what you did
- to be able to re-use the code and create variations
- for easy execution



```
1 ##### My first script #####
2 ##### October 2017 #####
3
4
5 # list workspace
6 ls()
7
8 # Reset R's brain
9 rm(list=ls())
10
11 # check wd
12 getwd()
13
14 #set wd
15 setwd("~/Users/dmarek/EducationSIB/Courses_2016/First_Steps_R_June2016/R_intro_course")
16
17 # confirm wd
18 getwd()
19
20 #load packages if needed (to do every time you launch your R session)
21 #library("boot")
22 #library("lattice")
23
24 #####
25 #####
26 #####
```

Notice the syntax highlighting

Create a new script and type code

- Create a new script using **File > New File > R script**. **Don't forget to save your script often.**
- By default, scripts are saved to the working directory.
- Files can be saved to other locations (**File -> Save As...**)
- Start **Typing code** at the top of the script

My first command:

2 + 3

- **Notice the syntax highlighting**
- **Comments** : “#” at the beginning of a line or before a command: helping text ; everything that follows is ignored by the during executing ; R does not support multi-line comments

Send Code From a Script to the Console

Run **individual lines**, one by one:

- In RStudio: put the cursor anywhere in a line, hit


Ctrl + enter (Windows)

Cmd + return (Mac)

or click the "Run" button

Tip: Run **part of a line** or **multiple lines**: **Highlight** the code, then proceed as above

Save, close and open scripts

- **Save a script:** File > Save or 
- **Close and open a script:** File > Close and File > Open File

Tips:

- Most of your code should be developed and saved in scripts.
- You can execute individual lines of code interactively while you are writing it.
- You can run the entire script once it is ready and debugged.


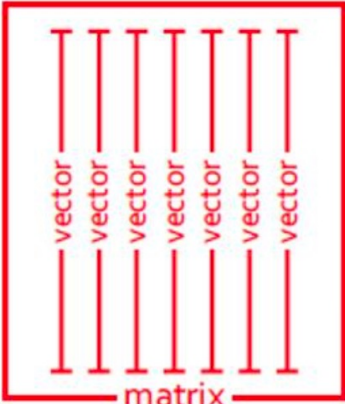
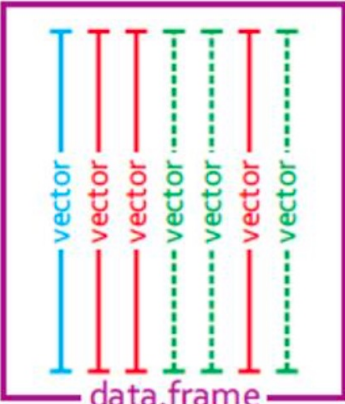
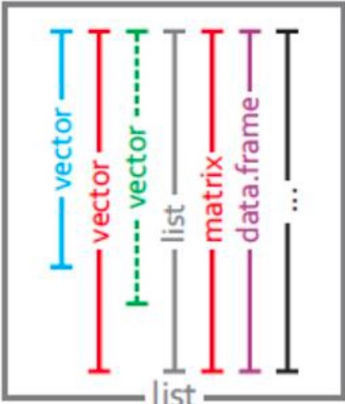
We continue working on the provided script

Download and open it

02

Syntax, data types and structures, importing data

Common object classes

	vector	matrix	data.frame	list
				
dimension	1	n	2	1
element data type	single	single	multiple	multiple
element data structure	atomic	atomic	vector	any
subsetting	x[i] x["name"] x[1:3]	x[i,j,...] x["row","col",...] x[,1:3,...]	x[i,j] x["row","col"] x\$colname	x[[i]] x\$colname

Example of a well-formated dataset

	A	B	C	D
1	Sample_ID	Age	Sex	Disease
2	M417	71	male	Healthy
3	M244	73	female	Tumor
4	M255	60	male	Healthy
5	M229	75	male	Tumor
6	M420	68	female	Healthy
7	M368	73	male	Healthy
8	M403	68	male	Tumor
9	M230	56	male	Tumor
10	M370	84	male	Tumor
11	M406	69	male	Tumor
12	M245	70	male	Tumor
13	M409	NA	female	Tumor
14	M395AR_dm	67	male	Tumor
15	PB	57	male	Healthy
16	M318	62	male	Healthy
17	M423	72	female	Tumor
18	M398_DMOS	61	female	Tumor
19	M233	74	male	Tumor
20	M381	57	male	Healthy
21	M408	65	male	Tumor
22	M402	68	male	Healthy

- A header line with variable names (4 variables, 1 in each column)
- No blank spaces in variable names (use _ instead)
- Variable names do not contain symbols other than _
- One observation per row
- No comments or other content around the data table
- Indicate missing values with NA

Example of a spreadsheet in Excel

Additional learning and practicing

Wandrille Duchemin's First Steps with R in Life Sciences (2 days):

It includes more on statistics!

<https://github.com/sib-swiss/first-steps-with-R-training/tree/master>

Introduction to statistics with R (3 days), for R beginners also:

<https://sib-swiss.github.io/Introduction-to-statistics-with-R/day1/>

Introduction to R for Cancer Scientists

<https://bioinformatics-core-shared-training.github.io/r-intro/index.html>

Glitr.org



Git repositories with
bioinformatics training material