

# Introduction to R for Life Sciences

**João Lourenço, Tania Wyss & Nadine Fournier**

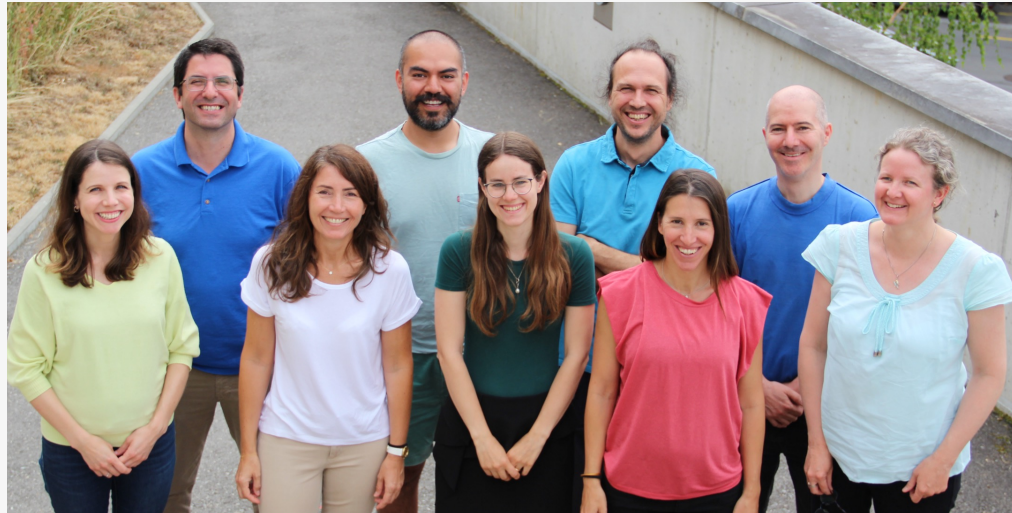
Translational Data Science – Facility

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](mailto: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

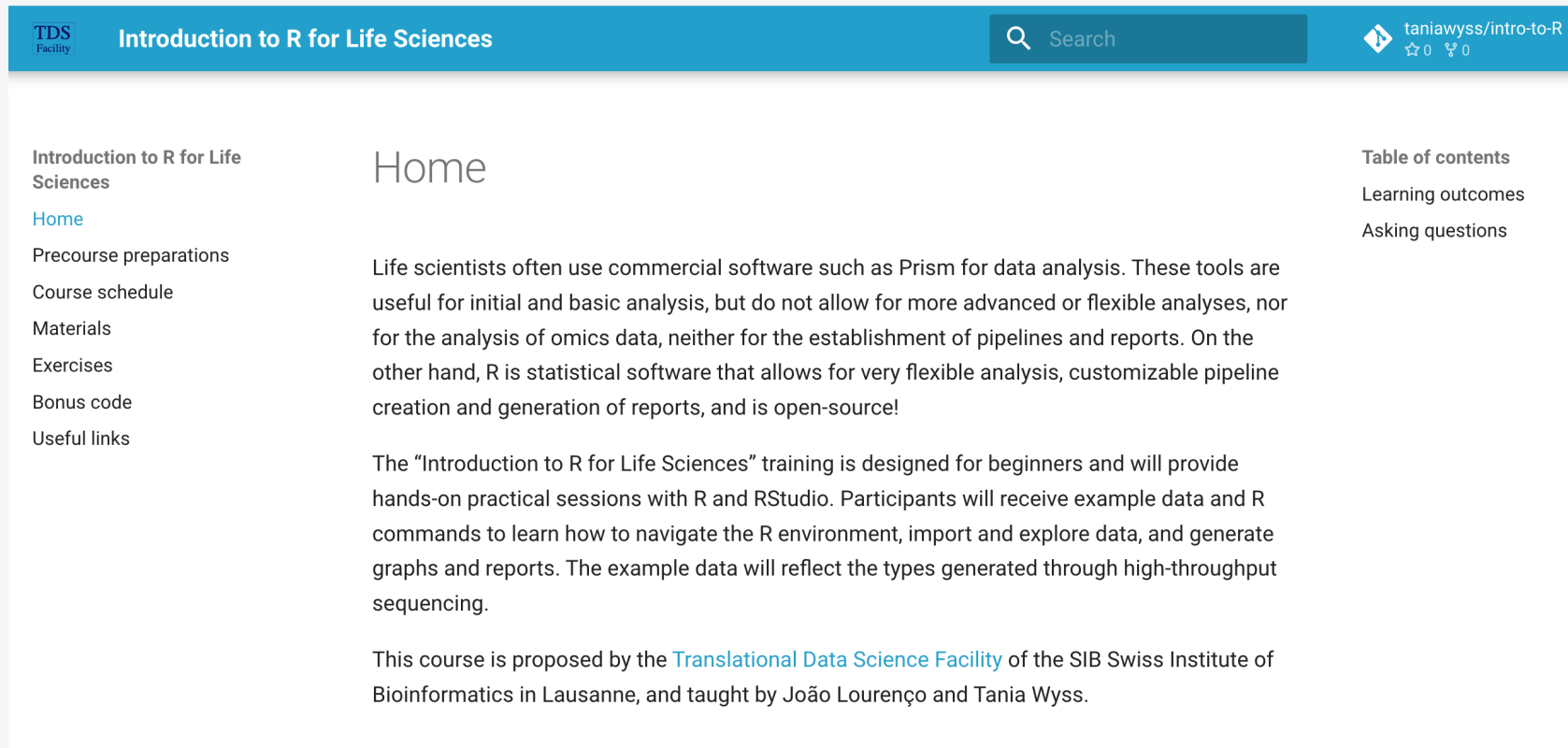


Photo by Scott Graham, Unsplash

# Course material

## 1. Website

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



The screenshot shows the homepage of the 'Introduction to R for Life Sciences' website. The header is blue and contains the TDS Facility logo, the title 'Introduction to R for Life Sciences', a search bar, and the repository name 'taniawyss/intro-to-R' with star and fork counts. The main content area is white and features a left sidebar with navigation links, a central 'Home' section with introductory text, and a right sidebar with a 'Table of contents'.

**Introduction to R for Life Sciences**

**Home**

Life scientists often use commercial software such as Prism for data analysis. These tools are useful for initial and basic analysis, but do not allow for more advanced or flexible analyses, nor for the analysis of omics data, neither 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, and is open-source!

The “Introduction to R for Life Sciences” training is designed for beginners and will provide hands-on practical sessions with R and RStudio. Participants will receive example data and R commands to learn how to navigate the R environment, import and explore data, and generate graphs and reports. The example data will reflect the types generated through high-throughput sequencing.

This course is proposed by the [Translational Data Science Facility](#) of the SIB Swiss Institute of Bioinformatics in Lausanne, and taught by João Lourenço and Tania Wyss.

**Table of contents**

- Learning outcomes
- 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.

A screenshot of the RGui (64-bit) window. The window title is "RGui (64-bit)". The menu bar includes "File", "Edit", "View", "Misc", "Packages", "Windows", and "Help". Below the menu bar is a toolbar with icons for file operations. The main area is the "R Console" window, which displays the following text:

```
R version 3.6.1 (2019-07-05) -- "Action of the Toes"  
Copyright (C) 2019 The R Foundation for Statistical Computing  
Platform: x86_64-w64-mingw32/x64 (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.  
  
> |
```

## 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 displays the RStudio interface with several panels and components:

- Editor (scripts):** The top-left panel shows an R script with the following code:

```
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 #####
```
- Environment and History:** The top-right panel shows the Global Environment with the following data objects:

Object	Class	Attributes
mice_data	Data Frame	50 obs. of 3 variables
mice_weight_HFD	Data Frame	29 obs. of 3 variables
mids	Matrix	num [1:2, 1] 0.7 1.9

Below the data objects, the Values section shows a list of variables and their values:

Variable	Class	Values
mean_weight_diet	num	[1:2(1d)] 28.7 37.1
mean_weight_genotype	num	[1:2(1d)] 33.7 33.4
n_weight_diet	int	[1:2(1d)] 21 29
n_weight_genotype	int	[1:2(1d)] 24 26
sd_weight_diet	num	[1:2(1d)] 2.61 5
sd_weight_genotype	num	[1:2(1d)] 4.69 6.92
- Console (or terminal):** The bottom-left panel shows the R console output:

```
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 panel shows the Files view of the current project directory: `~/TrainingSIB/Courses_2017/First_Steps_R_Oct2017/course_datasets`. It lists several CSV files with their sizes and modification dates:

Name	Size	Modified
..		
class.csv	402 B	Mar 6, 2016, 6:33 PM
etubiol.csv	2.4 KB	Mar 7, 2016, 12:13 PM
mammals_survey.csv	1.2 MB	Jan 30, 2017, 7:41 AM
melanoma_data.txt	5.9 KB	Jan 30, 2017, 7:41 AM
mice_data.csv	632 B	Feb 28, 2016, 4:15 PM
my_data_frame.csv	262 B	Mar 2, 2015, 4:49 PM
pigs.csv	206 B	Mar 6, 2016, 6:32 PM
smoker.csv	4.6 KB	Feb 26, 2015, 12:23 PM
snp.csv	618 B	Mar 6, 2016, 9:19 PM

# Console: The Command Line

~/TrainingSIB/Courses\_2017/First\_Steps\_R\_Oct2017/ ↗

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 ">" 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)
```

The (not always helpful) help pages:

```
> ?p.adjust
```

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.

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



**This causes R to execute it.**

02

## Working with script files

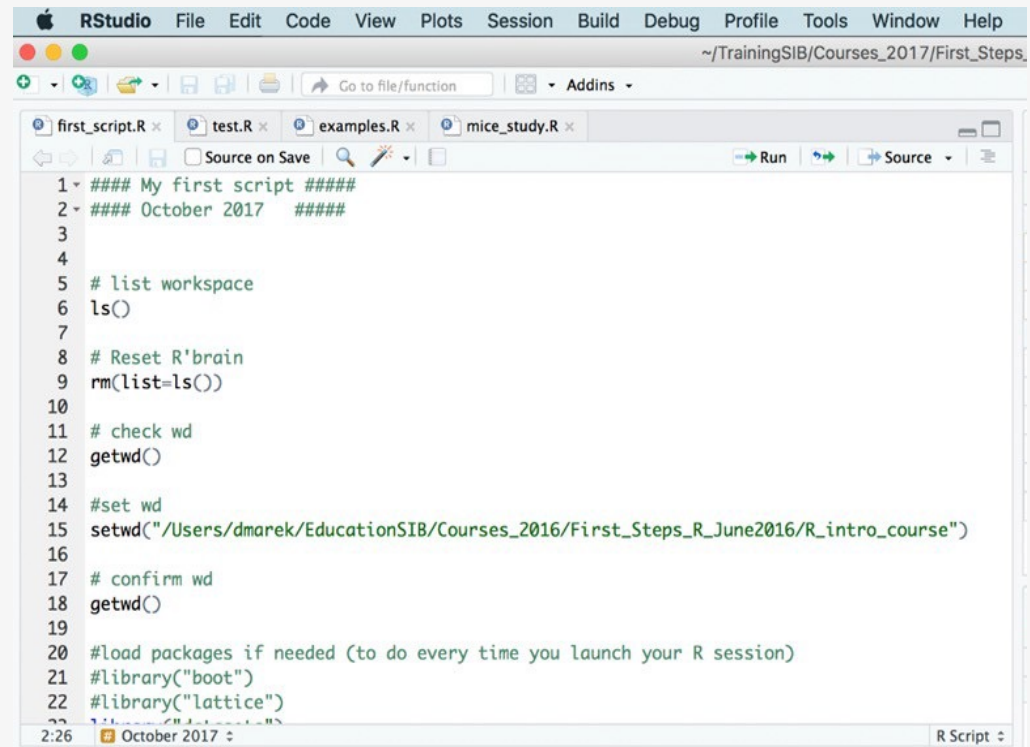


# Editor: Write code to a script file

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'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 ##### October 2017 #####
```

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.

# Let's work with the provided script!

Download it from the bottom of the Exercise page, and open it in R

 Exercises 🔍 Search

**Introduction to R for Life Sciences**

- Home
- Precourse preparations
- Course schedule
- Materials
- Exercises**
- Bonus code
- Useful links

```
sessionInfo()  
  
# Print the version of a specific package:  
packageVersion("stringi")
```

## Let's practice - Follow our script !

Download and open the provided commented script within R, run the commands and view the the output!

In our script, we included some sections with “fill in the blanks” exercises. You can either download the script without the solutions (first button), or the script with the solutions.

[Download script without solutions](#)

[Download script with solutions](#)

**⚠ Warning**


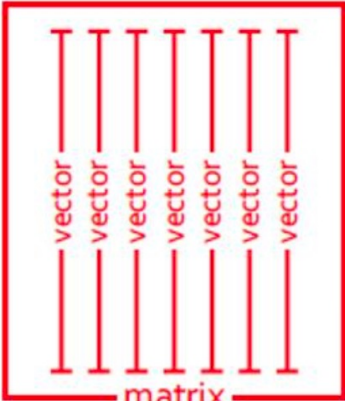
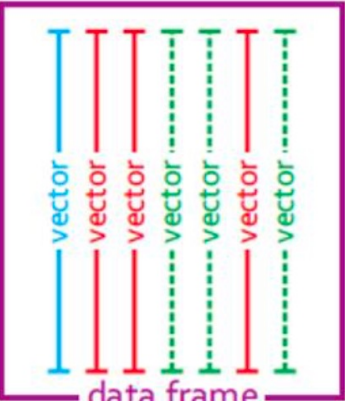
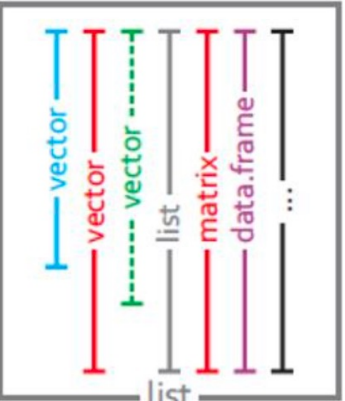
Make sure you have downloaded the csv files we import in the scripts from the [Materials](#) section.

**End of your first day with R, good job!**

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



03

## Graphics

# R graphics

R is powerful for plotting graphs and figures. Several plotting systems, including:

- base (i.e. graphics package, widely used, comes with basic R installation)
- ggplot2 (widely used in omics data analysis and others, implements the Grammar of Graphics, Wilkinson, Springer 2005)

*They have very different syntaxes, cannot be mixed, and need to be learned separately.*

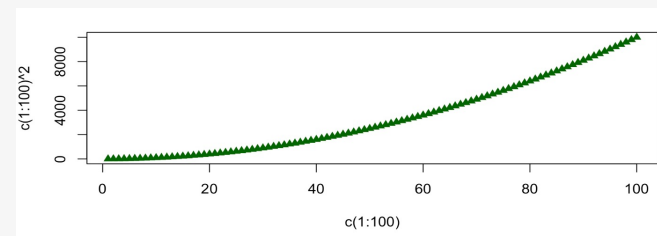
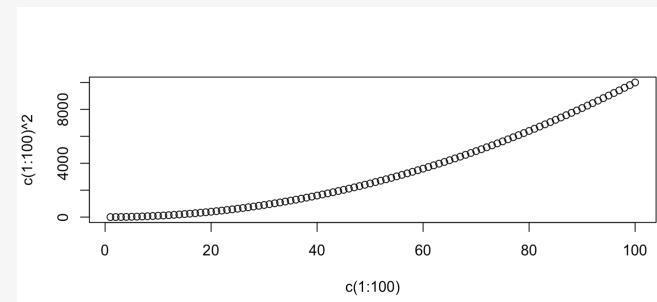
# R base plotting system

Plots are built up step by step with multiple function calls.

## High-level graphics functions:

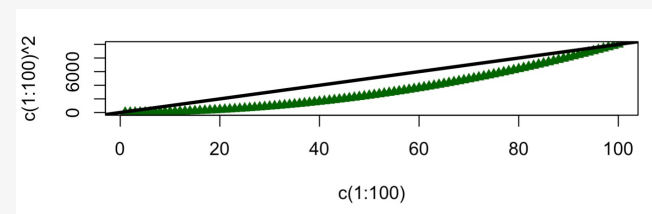
- Draw a new plot.  
> `plot(x=c(1:100), y=c(1:100)^2)`
- Tailor its appearance with optional arguments.

```
> plot(x=c(1:100), y=c(1:100)^2,  
      col="darkgreen", pch=17)
```



**Low-level graphics functions:** add graphical elements to an existing plot, piece by piece.

```
> plot(x=c(1:100), y=c(1:100)^2, col="darkgreen", pch=17)  
> abline(a=0, b=100, lwd=3)
```



# R colors (col)



























657 built-in color names  
Here is a subset -->

white	aliceblue	antiquewhite	antiquewhite1	antiquewhite2
antiquewhite3	antiquewhite4	aquamarine	aquamarine1	aquamarine2
aquamarine3	aquamarine4	azure	azure1	azure2
azure3	azure4	beige	bisque	bisque1
bisque2	bisque3	bisque4		blanchedalmond
blue	blue1	blue2	blue3	blue4
blueviolet	brown	brown1	brown2	brown3
brown4	burlywood	burlywood1	burlywood2	burlywood3
burlywood4	cadetblue	cadetblue1	cadetblue2	cadetblue3
cadetblue4	chartreuse	chartreuse1	chartreuse2	chartreuse3
chartreuse4	chocolate	chocolate1	chocolate2	chocolate3
chocolate4	coral	coral1	coral2	coral3
coral4	cornflowerblue	cornsilk	cornsilk1	cornsilk2
cornsilk3	cornsilk4	cyan	cyan1	cyan2
cyan3	cyan4	darkblue	darkcyan	darkgoldenrod
darkgoldenrod1	darkgoldenrod2	darkgoldenrod3	darkgoldenrod4	darkgray
darkgreen	darkgrey	darkkhaki	darkmagenta	darkolivegreen
darkolivegreen1	darkolivegreen2	darkolivegreen3	darkolivegreen4	darkorange
darkorange1	darkorange2	darkorange3	darkorange4	darkorchid
darkorchid1	darkorchid2	darkorchid3	darkorchid4	darkred
darksalmon	darkseagreen	darkseagreen1	darkseagreen2	darkseagreen3
darkseagreen4	darkslateblue	darkslategray	darkslategray1	darkslategray2
darkslategray3	darkslategray4	darkslategrey	darkturquoise	darkviolet
deeppink	deeppink1	deeppink2	deeppink3	deeppink4
deepskyblue	deepskyblue1	deepskyblue2	deepskyblue3	deepskyblue4

<https://www.nceas.ucsb.edu/sites/default/files/2020-04/colorPaletteCheatsheet.pdf>

<http://www.stat.columbia.edu/~tzheng/files/Rcolor.pdf>

# R plotting characters (pch)

<b>0</b> 	<b>1</b> 	<b>2</b> 	<b>3</b> 	<b>4</b> 	
<b>5</b> 	<b>6</b> 	<b>7</b> 	<b>8</b> 	<b>9</b> 	
<b>10</b> 	<b>11</b> 	<b>12</b> 	<b>13</b> 	<b>14</b> 	
<b>15</b> 	<b>16</b> 	<b>17</b> 	<b>18</b> 	<b>19</b> 	
<b>20</b> 	<b>21</b> 	<b>22</b> 	<b>23</b> 	<b>24</b> 	<b>25</b> 

## R line types (lty)



lty=1 or 'solid'



lty=2 or 'dashed'



lty=3 or 'dotted'



lty=4 or 'dotdash'



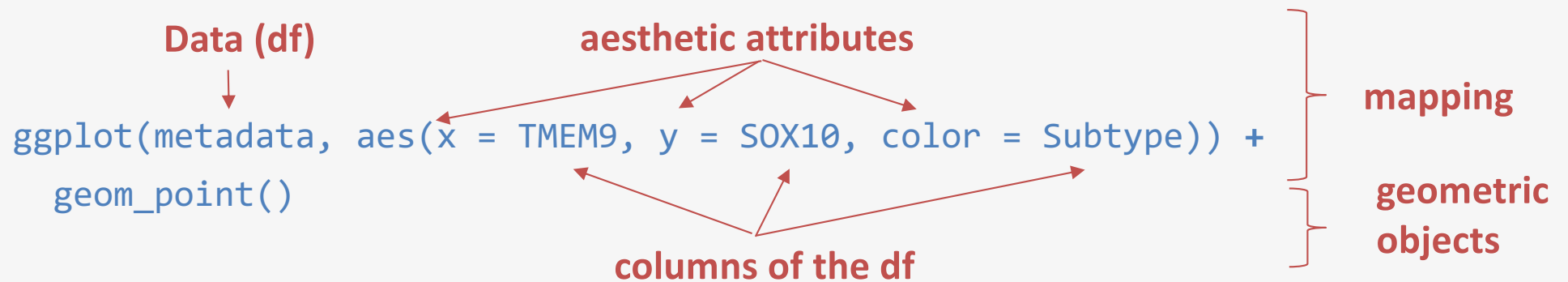
lty=5 or 'longdash'



lty=6 or 'twodash'

# ggplot2

An implementation of the **Grammar of graphics**: a graphics is a **mapping** from **data** to **aesthetic attributes** (coordinates, colors, shapes, sizes...) of **geometric objects** (points, lines, bars ...)



**Data should be in a data frame !**

# ggplot2

- The syntax (grammar) is very different from base R plotting functions.
- It builds a plot by adding layers of functions using the + sign
- The basic ggplot2 functions specify the data frame, the x,y coordinates, and the type of plot:

```
ggplot(dataframe, aes(x, y)) +  
geom_type()
```

Additional layers for full customizations  
are then added:

```
ggplot(dataframe, aes(x, y, color=factor)) +  
geom_type() +  
additional_layers()
```

**discrete x , continuous y**  
f <- ggplot(mpg, aes(class, hwy))



f + **geom\_col()**, x, y, alpha, color, fill, group, linetype, size



f + **geom\_boxplot()**, x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight



f + **geom\_dotplot**(binaxis = "y", stackdir = "center"), x, y, alpha, color, fill, group



f + **geom\_violin**(scale = "area"), x, y, alpha, color, fill, group, linetype, size, weight

## COLOR AND FILL SCALES (CONTINUOUS)

- o <- c + geom\_dotplot(aes(fill = x))
- o + **scale\_fill\_distiller**(palette = "Blues")
- o + **scale\_fill\_gradient**(low="red", high="yellow")
- o + **scale\_fill\_gradient2**(low = "red", high = "blue", mid = "white", midpoint = 25)
- o + **scale\_fill\_gradientn**(colors = topo.colors(6))  
Also: rainbow(), heat.colors(), terrain.colors(), cm.colors(), RColorBrewer::brewer.pal()

## SHAPE AND SIZE SCALES

- p <- e + geom\_point(aes(shape = fl, size = cyl))
- p + **scale\_shape()** + **scale\_size()**
  - p + **scale\_shape\_manual**(values = c(3:7))
  - o 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
  - p + **scale\_radius**(range = c(1,6))
  - p + **scale\_size\_area**(max\_size = 6)



# ggplot2

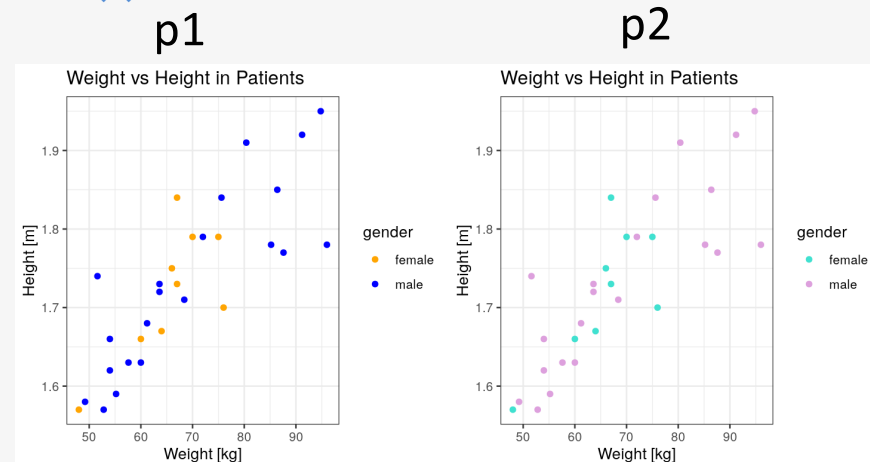
Many other packages offer additional ggplot2 functionalities

- Cowplot: Multi-panel figures; save the plot to an object, then display

```
p1 <- ggplot(dataset, aes(x, y, color=fact)) +  
  geom_type() +  
  additional_layers()
```

```
p2 <- ggplot(dataset, aes(x, y, color=fact)) +  
  geom_type() +  
  additional_layers()
```

```
install.packages("cowplot")  
library(cowplot)  
plot_grid(p1, p2, nrow=1)
```



- ggpubr: publication-ready plot customization, e.g. add T-test result on plot
- ggrepel: to avoid overlapping data point labels, e.g. in volcano plots
- ...

04

## Basic Statistics in R

# Statistical hypothesis testing

Two hypotheses in competition:

- H0: the NULL hypothesis (usually the most conservative – e.g., “no difference”)
- H1: the alternative hypothesis (usually the one we are actually interested in)

Example:

- H0: « There is no difference in the expression of a gene between two given subtypes of melanoma»
- H1: « The average expression of a gene is different in two given subtypes of melanoma»

Statistical test:

- Calculate test statistic
- Calculate associated p-value
- Check if p-value is small enough to reject H0, according to pre-defined significance level

# t-test

## Goal:

- Compare a continuous measure between two groups
- Is the difference between the two group means statistically significant?

## Assumptions:

- Observations are independent
- The two groups follow a normal distribution
- Homogeneity of variances (R uses Welch's t-test, which does not assume equal variance)

# Correlation

Measures the strength and direction of the relationship between two variables.

Correlation Coefficient ( $r$ ):

- Ranges from -1 to +1
- Direction
  - Positive correlation ( $r > 0$ ): As one variable increases, the other tends to increase
  - Negative correlation ( $r < 0$ ): As one variable increases, the other tends to decrease
- Strength
  - Perfect correlation:  $|r| = 1$
  - Strong correlation:  $0.7 < |r| < 1$
  - Moderate correlation:  $0.3 < |r| \leq 0.7$
  - Weak correlation:  $0 < |r| \leq 0.3$
  - No correlation (no consistent relationship between variables):  $r = 0$

**REMEMBER: Correlation does not imply causation !**

# Linear regression

Statistical method used to model the relationship between a dependent variable (Y) and one (or more) independent variable(s) (X).

Line of Best Fit:  $Y = \beta_0 + \beta_1 X$

- $\beta_0$ : Y-intercept (value of Y when  $X = 0$ )
- $\beta_1$ : Slope (change in Y for a one-unit increase in X)

Least Squares Method: Minimizes the sum of squared residuals

R-squared ( $R^2$ ):

- Measures how well the model fits the data
- Ranges from 0 to 1
- Higher values indicate better fit

# 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



# How to get data for practicing and playing

R contains many practice data sets (data frames), great for trying out functions.

## Display names of available data sets

```
> data(package = .packages(all.available = TRUE)) # lists  
data set names available in all installed packages
```

## Load and use a data set

```
> data(iris) # load the iris data (overwrite existing variable)  
> ?iris      # get information about the iris data  
> head(iris) # display top few lines of the iris data frame
```

	Sepal.Length	Sepal.width	Petal.Length	Petal.width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.2	setosa
6	5.4	3.9	1.7	0.4	setosa



# How to get data for practicing and playing

R can easily simulate data drawn from a given distribution. The function `rnorm()` generates normally distributed data.

Example:

```
>rnorm(10) #numeric vector with 10 values #drawn
           from normal distribution, #mean=0, sd=1
           (function defaults)
[1]  1.1053564  0.7937635  0.2743762  0.3574477 -0.7677099
[2]  0.5838973  0.6616164  0.1203090 -0.4060265  0.2778585
```

```
>rnorm(10, mean=10, sd=2) #customized mean and sd
[1]  6.253392  9.527140  9.398857 11.932284 11.472909
[2] 10.714245  7.656026 11.302829  9.332930 10.264157
```

If you want data from other distributions than normal:

`rpois()` for poisson, `rbinom()` for binomial (see R help)

sd: standard deviation