

SDA 2023 — Assignment 1

Solve the exercises below in RStudio (or some other IDE, but the assistants are most familiar with RStudio). RStudio is a graphical shell over R. Both programs are freeware, and can be installed on your own computer. **See the SDA Canvas page for useful links, including manuals in English;** the “particular R-manual” was recently translated for this course.

This assignment consists of 6 exercises. The aim of Exercises 1.2 and 1.3 is to get introduced to RStudio. If you have experience with R and RStudio, you may skip these introductory exercises.

Start solving this (and other) exercise(s) **well before** the practical classes! If you do, you will be able to ask your teaching assistant *relevant* questions instead of *trivial* ones that can be answered by a brief look into the Syllabus, lectures, or an R-manual or by some small preparational work.

Hand in Exercises 1.4–1.6. Try to solve these as efficiently as possible.

Make a concise report of your answers in *one single PDF file*, with only *relevant* R code in an *appendix*, i.e. the code that is needed to reproduce your findings. It is important to make clear in your answers how you have solved the questions. Graphs should look neat (label the axes, give titles, use correct dimensions etc.). Multiple graphs can be put into one figure using the command `par(mfrow=c(k,1))`, see `help(par)`.

Sometimes there might be additional information on what exactly should be handed in.

Carefully read the file AssignmentFormat.pdf on canvas.vu.nl.

Exercise 1.1 (Class) At the beginning of the practical class your teaching assistant will give a short demonstration of how to use RStudio. Be prepared to follow the same steps on your own laptop (make sure to have R and RStudio installed in advance!). Afterwards, continue on your own with Exercises 1.2 and 1.3 and complete every part of them.

Warning: In later exercises or assignments, your teaching assistant will not answer your syntax-related question if it is trivially resolved by carefully following Exercises 1.2 or 1.3. He/she will point you to these exercises in that case.

Exercise 1.2 Introduction to RStudio

Start up RStudio, and choose **File** → **New** → **R Script**. Now you should have 4 windows.

top left script window — typing, editing, saving, executing R-commands

bottom left console window — for direct typing and executing commands (**no saving!**)

top right workspace window — overview of known variables, import datasets

bottom right plot window — for graphics (view, save, etc) and help-function

The sign ‘>’ at the beginning of a line in the console window (bottom left) is the R prompt. It is followed by the commands that you type. In case a command is not yet finished, the next line will show a + sign, and you can continue your command after this sign.

Below you find a set of introductory R commands, that shows some of its possibilities. The right column contains some explanation of the corresponding command in the left column. Type these commands directly in the console window and see what you get.

> x = 1:20 (or x <- 1:20)	make a vector x with values 1, 2, 3, ..., 20.
> x	print value of x on the screen.
> m = matrix(x,4,5,byrow=T)	create a matrix m with 4 rows and 5 columns
	with the values of x ordered row-wise.
> m	print matrix m on the screen.
> m[2,3]	print element (2,3) of m on the screen.
> m[2,]	print all elements of 2 nd row of m .

<code>> m[,3]</code>	print all elements of 3 rd column of <code>m</code> .
<code>> y = sample(1:100,20)</code>	generate random sample of size 20 from the numbers 1, 2, 3, ..., 100.
<code>> z = x+y</code>	compute the sum of <code>x</code> and <code>y</code> coordinate-wise.
<code>> y = x+ 2*y</code>	transform <code>y</code> coordinate-wise.
<code>> cbind(x,y)</code>	form a matrix with columns <code>x</code> and <code>y</code> and print the result.
<code>> z <- c(NA, 1/0, 0/0)</code>	creates a vector of NA (not available), Inf (infinite), NaN (not a number)
<code>> is.na(z)</code>	check for missing values and print the result.
<code>> plot(x,y)</code>	plot <code>y</code> against <code>x</code> .
<code>> abline(100,2)</code>	add the line <code>y=100+2*x</code> to the last plot.

The drawback of typing directly in the console window is the lack of saving the typed commands. The script window (top left) is very useful if you want to type, edit (e.g. correct your typo's) and save code. You can execute lines in the script window by pressing **Ctrl+Enter** (on Mac: **Cmd+Enter**). Execute the remainder of this introductory R commands from the script window. Try the **File** → **Save** option, you will need it later on!

<code>> set.seed(1234)</code>	sets the seed for random data generation
<code>> x = rnorm(50,0,sqrt(2))</code>	generate random sample of size 50 from normal distribution with mean 0 and variance 2.
<code>> y = rnorm(50)</code>	idem from standard normal distribution.
<code>> mean(x)</code>	compute mean of the values in <code>x</code> .
<code>> sd(x)</code>	compute standard deviation of the values in <code>x</code> .
<code>> var(x)</code>	compute variance of the values in <code>x</code> .
<code>> cor(x,y)</code>	compute correlation between <code>x</code> and <code>y</code> .
<code>> x[x<0]</code>	select negative elements in <code>x</code> .
<code>> sum(x<0)</code>	count number of negative elements in <code>x</code> .
<code>> hist(x,prob=T)</code>	plot a scaled histogram.
<code>> help(hist)</code>	give documentation about the function <code>hist</code> .
<code>> ?hist</code>	the same.
<code>> f <- function(x){x*x}</code>	defines a function <code>f(x)</code> that does the same as <code>x^2</code>
<code>> u = seq(-5,5,0.1)</code>	form sequence of points between -5 and 5 with step size 0.1.
<code>> v = dnorm(u,0,sqrt(2))</code>	compute density of normal distribution with mean 0 and variance 2.
<code>> lines(u,v)</code>	add plot of computed normal density.
<code>> {hist(x,xlim=c(-6,6),prob=T)</code>	repeat plot of scaled histogram on larger interval, but do not yet execute the command.
<code>+ lines(u,v)}</code>	plus additional command, and execute both.
<code>> plot(sort(x), 1:50/50,</code>	plot empirical distribution function of <code>x</code> .
<code>+ type="s",ylim=c(0,1),</code>	
<code>+ xlab="x",ylab="the ecdf of x")</code>	
<code>> lines(u,pnorm(u,0,sqrt(2)))</code>	add true distribution function.
<code>> w = seq(-pi,pi,length=100)</code>	form regular grid of 100 points.
<code>> plot(cos(w),sin(w),type="l")</code>	draw a circle.

Navigate through your plots, using the arrow buttons in the top line of the plot window (bottom right). Use the **Export** button to save them as picture files.

Exercise 1.3 *Vectors and matrices*

- Use the function `c()` to create a vector `x` that consists of the numbers 23, 0.1, -5.15.
- Add to this vector the number -28, using again `c` or `append`.
- Order the elements of `x` from small to large using `sort` and call the vector of ordered numbers `y`.
- Multiply each element of `x` by 4.
- Round each element of `x` to one decimal place using `round`.
- Select all positive elements of `x`.
- Change the third element of `x` into 7.
- Create a vector `z` consisting of the sequence of numbers between 2 and 4 with step size 0.1. Create a matrix `m` with seven rows and three columns in which the 21 numbers in `z` are ordered columnwise.
- Extract the element on the second row and in the first column of `m`.
- Extract the third column of `m`.
- Compute the mean value of the numbers in `m`.
- Compute the mean value of each column of `m` by using the R function `apply`.

For the following exercises the R functions `hist`, `stem`, `boxplot`, `plot`, `summary`, `set.seed`, `rlnorm`, `rexp`, `list`, `quantile` and `apply` may be helpful. Use `help(yyyy)` for the manual of any function called `yyyy`.

Exercise 1.4 (**.RData file hand-in**) Write a function `norm(n,mu,sigma)` that

- sets the seed to 2023210,
- draws a random sample of size `n` from the normal distribution with parameters `mean=mu`, `sd=sigma`,
- creates a list `mylist` that contains

`quants` = a vector consisting of the 5%, 50%, and 95% quantiles of the drawn sample (use the function `quantile` with the default `type=7`),

`loc` = the sample mean of the drawn sample,

`spread` = the sample standard deviation of the drawn sample,

`stud_no` = a vector that contains the student numbers of you (and your group partner);

- stores the list in an .RData file (use `save(mylist, file="[PATH]/myfile1.RData")`, where `[PATH]` stands for the path on your computer where you wish to save the .RData file.

Call `norm` with the parameters `n=100`, `mu=1`, `sigma=1`.

Hand in: Submit your .RData file created by a call of the above-described function to the (dummy) assignment called “RData A1” on Canvas.

In addition, the R-code of your function must be in the appendix of the regular report you will submit to “Assignment 1” on Canvas.

Exercise 1.5

a. Write a function `CLT_unif(n,m)` that

- draws n samples $\mathbf{x}_1, \dots, \mathbf{x}_n$ of size m from the uniform distribution $U(0, 1)$ on the interval $(0, 1)$;
- plots a scaled histogram of the n sample-specific sample means (scaled to the probability level), i.e., first compute for each of the n samples the mean of m numbers, then plot the histogram;
- plots in the same figure the density of the normal distribution with mean 0.5 and standard deviation $1/\sqrt{12m}$ (as a red colored curve).

Hint: Use the function `lines` with the additional options `type="l", col="red"` to plot a function in red.

Use the functions `runif`, `hist` and `dnorm` (see e.g. `help(hist)` for back-up information). Use the function `lines` to combine two graphs in one figure. Furthermore, the graph should look appropriate: use a proper title and axis labels (you could use the function `paste` for this), make sure that nothing is cropped.

Note: Always make sure to use proper labels and captions for your plots!

b. Make plots for the combinations $(n, m) \in \{(50, 30), (50, 200), (300, 30), (300, 200)\}$. Indicate the influence of each of the two parameters n and m on how well the histogram resembles the normal density.

Hand in: the (final) code of your function (in the appendix!), your answer to the last question, and the 4 graphical realizations of the function (as explained above) such that they support your answer.

Exercise 1.6 We wish to make data summaries of the military expenses of states per citizen in 2020. In addition, we will investigate the relationship between such expenses in 2020 and those in 1988. To this end, we use the dataset `military-spending-per-capita.csv`¹ (to be found on Canvas).

Proceed as follows: read the data into your R environment, e.g., by using the R commands provided in the .R file `help_Assignment_1.R`.

a. Make univariate numerical and graphical summaries of the military expenses per capita in 2020.

Note: for the graphical summaries, create at least a histogram and a boxplot.

b. Make bivariate numerical and graphical summaries of the corresponding 1988 and 2020 data, only for those countries for which both numbers are available.

Note: you should always describe your findings in words!

¹source: <https://ourworldindata.org/>, accessed on Jan. 30, 2023