

## SDA 2021 — 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 4 exercises. The aim of Exercises 1.1 and 1.2 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 1.3 & 1.4.** Try to solve these (completely) 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 AssignmentFormat2021.pdf on canvas.vu.nl.**

### Exercise 1.1 *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 <b>x</b> with values 1, 2, 3, ..., 20.
> x	print value of <b>x</b> on the screen.
> m = matrix(x,4,5,byrow=T)	create a matrix <b>m</b> with 4 rows and 5 columns with the values of <b>x</b> ordered row-wise.
> m	print matrix <b>m</b> on the screen.
> m[2,3]	print element (2,3) of <b>m</b> on the screen.
> m[2,]	print all elements of 2 <sup>nd</sup> row of <b>m</b> .
> m[,3]	print all elements of 3 <sup>rd</sup> column of <b>m</b> .
> y = sample(1:100,20)	generate random sample of size 20 from the numbers 1, 2, 3, ..., 100.
> z = x+y	compute the sum of <b>x</b> and <b>y</b> coordinate-wise.
> y = x+ 2*y	transform <b>y</b> coordinate-wise.
> cbind(x,y)	form a matrix with columns <b>x</b> and <b>y</b>

<code>&gt; z &lt;- c(NA, 1/0, 0/0)</code>	and print the result. creates a vector of NA (not available), Inf (infinite), NaN (not a number)
<code>&gt; is.na(z)</code>	check for missing values and print the result.
<code>&gt; plot(x,y)</code>	plot y against x.
<code>&gt; abline(100,2)</code>	add the line $y=100+2*x$ 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>&gt; x = rnorm(50,0,sqrt(2))</code>	generate random sample of size 50 from normal distribution with mean 0 and variance 2.
<code>&gt; y = rnorm(50)</code>	idem from standard normal distribution.
<code>&gt; mean(x)</code>	compute mean of the values in <b>x</b> .
<code>&gt; sd(x)</code>	compute standard deviation of the values in <b>x</b> .
<code>&gt; var(x)</code>	compute variance of the values in <b>x</b> .
<code>&gt; cor(x,y)</code>	compute correlation between <b>x</b> and <b>y</b> .
<code>&gt; x[x&lt;0]</code>	select negative elements in <b>x</b> .
<code>&gt; sum(x&lt;0)</code>	count number of negative elements in <b>x</b> .
<code>&gt; hist(x,prob=T)</code>	plot a scaled histogram.
<code>&gt; help(hist)</code>	give documentation about the function <b>hist</b> .
<code>&gt; ?hist</code>	the same.
<code>&gt; f &lt;- function(x){x*x}</code>	defines a function <b>f(x)</b> that does the same as $x^2$
<code>&gt; u = seq(-5,5,0.1)</code>	form sequence of points between -5 and 5 with step size 0.1.
<code>&gt; v = dnorm(u,0,sqrt(2))</code>	compute density of normal distribution with mean 0 and variance 2.
<code>&gt; lines(u,v)</code>	add plot of computed normal density.
<code>&gt; {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>&gt; plot(c(-5,sort(x),5),</code>	plot empirical distribution function of <b>x</b> on (-5,5). <sup>1</sup>
<code>+ c(0,1:50/50,1),</code>	
<code>+ type="s",ylim=c(0,1),</code>	
<code>+ xlab="x",ylab="the ecdf of x")</code>	
<code>&gt; lines(u,pnorm(u,0,sqrt(2)))</code>	add true distribution function.
<code>&gt; w = seq(-pi,pi,length=100)</code>	form regular grid of 100 points.
<code>&gt; plot(cos(w),sin(w),type="l")</code>	draw a circle.

<sup>1</sup> Note: first and last value of vector to be plotted on x-axis need to be smaller than **min(x)** and larger than **max(x)**, respectively. You may need to replace the values -5 and/or 5 to account for this.

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 1.2 *Vectors and matrices*

- a. Use the function `c()` to create a vector `x` that consists of the numbers 23, 0.1, -5.15.
- b. Add to this vector the number -28, using again `c` or `append`.
- c. Order the elements of `x` from small to large using `sort` and call the vector of ordered numbers `y`.
- d. Multiply each element of `x` by 4.
- e. Round each element of `x` to one decimal place using `round`.
- f. Select all positive elements of `x`.
- g. Change the third element of `x` into 7.
- h. 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.
- i. Extract the element on the second row and in the first column of `m`.
- j. Extract the third column of `m`.
- k. Compute the mean value of the numbers in `m`.
- l. 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`, and `apply` may be helpful. Use `help(yyyy)` for the manual of any function called `yyyy`.

### Exercise 1.3

- a. Write a function `exp(n,rate)` that
  - draws a sample of size `n` from the exponential distribution with rate `rate`,
  - plots a scaled histogram of the sample (scaled to the probability level; with black colored lines and white filling, i.e. the default).
  - plots in the same figure the density of the `Exp(rate)`-distribution (as a red colored curve).

Use the functions `rexp`, `hist` and `dexp` (see `help(rexp)` for back-up information). Make sure that the area under the histogram equals 1, like the area under the density (see `help(hist)`). 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.

- b. Play around with different values of `n`  $\geq 10$  and `rate`  $\geq 3$ . Indicate the influence of each parameter on how well the histogram resembles the true density.

**Hand in:** the (final) code of your function, your answer to the last question, and 3 or 4 graphical realisations of the function with the corresponding arguments (`n,rate`) as illustration to your answer.

**Exercise 1.4** We wish to make data summaries of the new coronavirus cases (smoothed over 7 days) in Europe and also compare the numbers in December 2020 and January 2021. To this end, we are going to use the dataset `owid_covid_data.csv`<sup>1</sup> to be found on Canvas.

Proceed as follows: first read the data into your R environment, then extract the relevant information: the `new_cases_smoothed_per_million` for all countries in Europe for the dates December 20, 2020, and January 20, 2021; you might want to store the December and January values in two different vectors.

- a. Make univariate numerical and graphical summaries of each of the December and the January samples. Use these to compare both samples.

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

- b. Make bivariate numerical and graphical summaries of the samples.

*Note: you should always describe your findings in words!*

*Tip: the data can be loaded using the button **Import dataset** in the workspace window (top right). Alternatively, and perhaps preferably because you could make use of many important options, you could load the data via the function `read.table`; see e.g. Lecture 1 for details.*

---

<sup>1</sup>actual source: <https://ourworldindata.org/coronavirus-data>