

# *Applied Statistics (A) - Lecture 1*

## *R fundamentals and programming*

Gaia Bertarelli <sup>1</sup>    [gaia.bertarelli\(at\)santannapisa.it](mailto:gaia.bertarelli(at)santannapisa.it)

<sup>1</sup>Institute of Management  
Sant'Anna School of Advanced Studies

January 2023

# The presentation at a Glance

- 1 R fundamentals and programming
  - Basic tools
- 2 Exercises

# Why R

- ★ The software environment R is nowadays the most commonly used software in the statistical world,
- ★ R was founded by Ross Ihaka and Robert Gentleman in 1994/1995,
- ★ Since 1997, it has been internationally developed and distributed from Vienna over the Comprehensive R Archive Network (CRAN),
- ★ R is open-source (*free*),
- ★ R is object-oriented, which allows for efficient handling of complex structures (*flexible and object oriented*),
- ★ R provides a powerful interface to integrate programs written in other languages (*easily interfaceable*),
- ★ R has a huge, active and constantly evolving users-based community (*community-supported*),
- ★ R has amazing graphical capabilities (*highly graphically skilled*).

R, with its ever-growing number of **interfacing packages**, provides a great help in making integration easier, even without specific knowledge of advanced computer science.

# References about R

- ★ Official manuals here: <http://cran.r-project.org/>
- ★ Course materials (Hafner, 2019),
- ★ the whole *Use R!* series published by Springer, in which you can find specialized texts on variety of topics (we will use *Introducing Monte Carlo Methods with R* Authors: Robert, Casella.)
- ★ the funny guide *The R Inferno*, by Patrick Burns, available here: [http://www.burns-stat.com/pages/Tutor/R\\_inferno.pdf](http://www.burns-stat.com/pages/Tutor/R_inferno.pdf).  
'If you are using R and you think you're in hell, this is a map for you'.

# The RStudio IDE

- ★ When programming, the need to resort to other tools arises: for example, external text editors, file managers or web browsers.
- ★ **IDE** (Integrated Development Environment)
- ★ RStudio (<https://www.rstudio.com/>)
- ★ In just one screen:
  1. a powerful text editor with **syntax highlighting**, brackets matching (VERY useful) and buttons/keyboard combinations to submit code snippets to the console directly (i.e., no need for copy-and-paste nonsense)
  2. a **workspace browser**, a **data viewer** and the **commands history**
  3. a **file manager**, a **package manager**, a **plot** tab (where you can browse all the plot you created, back and forth) and an **integrated help tab**.

# Organize your work

- ★ There is no general rule on how to organize your working directories
- ★ I create a single directory for each project **ProjectName** and create then the following sub-directories
  - **ProjectName/Code**: this is where to save all the *.R* code snippets
  - **ProjectName/Data**: *.txt*, *.xls*, *.csv*, *.dat*, *.spss*, ..., raw and formatted data files in general go here
  - **ProjectName/Literature**: did you find a paper related to your project? Put it here!
  - **ProjectName/Out**: any output of your work, e.g. plots, text summaries, ...
  - **ProjectName/Varia**: anything else that doesn't belong to the other directories (project documentation, call for papers, to-do lists, ...).

# Help functions

R provides a wealth of functions, datasets and packages (more than 6000 now), definitely too many to be able to keep trace of every functionality, syntax or purpose. Some commands that might come in handy are:

- ★ `help()`: need to find out what function `xxx` does? simply enter `help(xxx)` in the command console (`?xxx` is an equivalent way to do so). Were you interested in the help for operators involving special symbols in their syntax (such as the matrix multiplication operator `%*%`), do not forget to put it between apostrophes: `?'%*%'`
- ★ `help.search()`: equivalently, `??`, allows for searching the help system for documentation matching a given character string. Say we want to find every function in every loaded package that deals with optimization, then we should enter `??optimization`; this command also accepts regular expressions

- ★ `apropos()`: accepts a string as an argument, and returns a character vector giving the names of all objects in the search list matching it; it might be useful to find, for example, all the functions containing the word `glm` (`apropos('glm')`)
- ★ `data()`: looking for some data? just enter `data()` on the command line and a list of all the datasets available in the default package datasets and in the currently loaded libraries, together with a short description thereof, will appear. These datasets can then be readily examined with the help command (`?nameofthedatast`), or loaded into the working space (`data(nameofthedatast)`).

If you are having coding problems

- ★ <http://stackoverflow.com>
- ★ <http://www.r-project.org/mail.html>.



# Basic Rules

- ★ R is **case sensitive**; A and a are different symbols and would refer to different variables.
- ★ Commands are separated either by a **semi-colon** (;), or by a **newline**
- ★ Comments can be put almost anywhere, starting with a **hashmark** (#), everything to the end of the line is a comment.
- ★ If a command is not complete at the end of a line, R will give the symbol "+"
- ★ R is an **object-oriented** program: every operation is made on and produces objects.
- ★ All objects in R have a class, reported by the function **class()**. Possible class are: numeric, logical, character, list, matrix, array, factor and data.frame.

# Basic Rules

- ★ Assignment is performed by the "<-" or the "=" operator.
- ★ Names defining objects cannot contain spaces or mathematical operators/special characters (except for the dot .), nor can they begin with a number;
- ★ Some peculiar values:
  - NA (Not Available) is the code denoting a missing numerical or character element (warning: "NA" is a valid character string);
  - NaN (Not a Number) is the result of impossible or undefined expressions like a division by zero;
  - Inf and -Inf denote infinity
- ★ These and other R keywords (for,while,if,TRUE,FALSE etc.) are not available for assignment.

- ★ A **vector** is a sequence of data elements of the same basic type.
- ★ To define a vector we use the function **c()**
- ★ Basic mathematical functions can also be applied to vectors. Such functions are performed **element-by-element**, i.e. elementwise
- ★ In other cases - like for basic statistical summary measures - the whole vector is the input of the function.

# Factors

- ★ are objects encoding **categorical variables**;
- ★ are often **used to group other** (usually quantitative) variables;
- ★ are defined by the command **factor()**;
- ★ have a **reference level** which can be changed.

# Lists

- ★ are set of objects with **different nature/dimension**;
- ★ are defined by the command **list()**;
- ★ are useful to summarise an analysis or as output of **complex function**.

# Matrices and Arrays

- ★ A matrix is a collection of data elements arranged in a two-dimensional rectangular layout:
- ★ a matrix is defined with the function `matrix()`
- ★ Arrays are the R data objects which can store data in more than two dimensions
- ★ An array is defined with the function `array()`.

# Dataframes

- ★ Data frame is a two dimensional data structure in R. It is a **special case of a list which has each component of equal length**.
- ★ Data frame is of particular importance in data analysis. It represents the matrix of data where **each row is an observation and each column is a variable** (the variables may be of different type);
- ★ Dataframes are defined by command: **data.frame()**.

# R libraries

- ★ All R functions and datasets are stored in packages (also called **libraries**).  
Only when a package is loaded are its contents available.
- ★ A package is a collection of **previously programmed functions**.
- ★ There are about 25 packages supplied with R (called "standard" and "recommended" packages)
- ★ The others are available through the CRAN family



# Packages handling

- ★ R basic installation comes with a number of packages that have been deemed essential.
- ★ However, you will often need to resort to specific functions and routines, that can be downloaded by one of the many repositories scattered around the world (wide web).
  - `install.packages()`: as the name states, it allows to install new packages, provided the exact name is matched, as is `install.packages('nameofthepackage')`
  - `library()`: if the package you need is available on your computer, you have to load it in the current Rsession, in order for its contents to be at your disposal; entering `library(nameofthepackage)` will do exactly this, making all the functions and datasets in the package available for usage. The `require()` command accomplishes the same task, but is intended specifically for use inside functions.

# Objects and indexing

- ★ R is an object-oriented programming (OOP) language
- ★ An R object is basically a category, examples are: scalar, vector, matrix, array, list, time series, data frames, functions and graphics
- ★ An object, say `x`, is characterized by:
  - a `mode(x)` that describes its contents
  - and a `class(x)` that describes its nature: class is a property assigned to an object that determines how generic functions operate with it.
- ★ An object's `mode` might be one of the following:
  - `null`: empty object
  - `logical`: TRUE or FALSE
  - `numeric`: a real number, such as  $\pi$ , 2 or  $\sqrt{5} - 1$
  - `complex`: a complex number, such as  $2 + 5i$
  - `character`: a string, such as 'normal', 'green', 'M' or 'y=a+bx'

## R code

```
1 > x <- 10
2 > x
3 [1] 10
4 > mode(x)
5 [1] "numeric"
6 > x <- 10>5
7 > x
8 [1] TRUE
9 > mode(x)
10 [1] "logical"
```

## R code

```
1 > x<-1:16
2 > mode(x)
3 [1] "numeric"
4 > class(x)
5 [1] "integer"
6 > dim(x)<-c(4,4)
7 > mode(x)
8 [1] "numeric"
9 > class(x)
10 [1] "matrix"
11 > is.numeric(x)
12 [1] TRUE
13 > mode(x) <- "character"
14 > mode(x)
15 [1] "character"
16 > class(x)
17 [1] "matrix"
18 > x <- factor(x)
19 > class(x)
20 [1] "factor"
21 > mode(x)
22 [1] "numeric"
```

Many functions exist, that can create objects of a specific kind, for example:

- ★ `c()` defines a row vector of given elements
- ★ `numeric()` defines an empty vector of given length
- ★ `matrix()` defines a matrix of given number of rows and columns, and given elements
- ★ `array()` defines an array of given dimensions and elements
- ★ `data.frame()` defines a data-frame object
- ★ `list()` defines a list object, i.e. a special array that can include elements with various modes

The function `str(x)` will compactly display the internal **structure** of an R object, and is a very useful command, especially when working with complex ones.

## R code

```
1 > y <- list(a=c(1,2,3),b=matrix(c(1,2,3,4,5,6),3,2))
2 > y
3 $a
4 [1] 1 2 3
5
6 $b
7      [,1] [,2]
8 [1,]    1    4
9 [2,]    2    5
10 [3,]    3    6
11
12 > mode(y)
13 [1] "list"
14 > str(y)
15 List of 2
16 $ a: num [1:3] 1 2 3
17 $ b: num [1:3, 1:2] 1 2 3 4 5 6
```

## R code

```
1 > y$a           #extracts component 'a' from list 'y'
2 [1] 1 2 3
3 > y$b           #extracts component 'b' from list 'y'
4      [,1] [,2]
5 [1,]    1    4
6 [2,]    2    5
7 [3,]    3    6
8 > dim(y$b)      #what are the dimensions of matrix 'b'?
9 [1] 3 2         #3 rows and 2 columns
10 > y$b[1, ]      #extracts the first row of matrix 'b'
11 [1] 1 4
12 > y$b[, 2]      #extracts the second column of matrix 'b'
13 [1] 4 5 6
14 > y$b[1:2, ]    #extracts rows from 1 to (:) 2 of matrix 'b'
15      [,1] [,2]
16 [1,]    1    4
17 [2,]    2    5
18 > y$a[3]        #extracts the third element of vector 'a'
19 [1] 3
```

It is possible to assign values to specific elements of an object:

## R code

```
1 > y$a[3] <- 10 #assign to the third element of vector 'a'
2 > y$a          #the value 10
3 [1] 1 2 10
4 > y$b[,1] <- c(2,3,7) #replace the first column of matrix 'b'
5 > y$b          #with a vector having components (2, 3, 7)
6               [,1] [,2]
7 [1,]          2     4
8 [2,]          3     5
9 [3,]          7     6
```



An array is a matrix in dimensions higher than 2 (a *tensor*). Indexing works in the same way, allowing to extract and manipulate information at every level by means of the square brackets operators `[,...,]`.

```
1 > z <- array(runif(24),dim=c(4,3,2))
2 > #creates an array composed by two 4x3 matrices
3 > #filled by realizations from a U(0,1) r.v.
4 > z
5 , , 1
6           [,1]      [,2]      [,3]
7 [1,] 0.50948231 0.09637384 0.3511716
8 [2,] 0.76222931 0.49445906 0.9239261
9 [3,] 0.07748602 0.24500671 0.9360242
10 [4,] 0.14746812 0.63679512 0.1754347
11
12 , , 2
13           [,1]      [,2]      [,3]
14 [1,] 0.4685381 0.9641261 0.9762254
15 [2,] 0.3219655 0.5849373 0.7358713
16 [3,] 0.5257197 0.8853259 0.5509927
17 [4,] 0.7837449 0.4370338 0.9350366
```

## R code

```
1 > z[, ,1] #extracts the first matrix
2           [,1]      [,2]      [,3]
3 [1,] 0.50948231 0.09637384 0.3511716
4 [2,] 0.76222931 0.49445906 0.9239261
5 [3,] 0.07748602 0.24500671 0.9360242
6 [4,] 0.14746812 0.63679512 0.1754347
7 > z[1, ,2] #extracts the first row of the second matrix
8 [1] 0.4685381 0.9641261 0.9762254
9 > z[,2,1] #extracts the second column of the first matrix
10 [1] 0.09637384 0.49445906 0.24500671 0.63679512
11 > z[,1,1] <- 1:4
12 > #assign the values from 1 to 4 to the first column of
13 > #the first matrix
```

- ★ R math operators usually work element-wise (to perform an operation on each element of a vector while doing a computation),
- ★ Element-wise operations are more intuitive than vectorwise operations, because the elements of one matrix map clearly onto the other, and to obtain the result, you have to perform just one arithmetical operation.

```
1 > x <- c(1,2,5,8)           #defines a vector
2 > 2*x                       #multiplies every element of x by 2
3 [1]  2  4 10 16
4 > x^2                       #squares every element of x
5 [1]  1  4 25 64
6 > sqrt(x)                   #extracts square roots
7 [1] 1.000000 1.414214 2.236068 2.828427
```

Indexing allows to apply specific operations on chosen parts of an object

## R code

```
1 > x <- matrix(1:25,5,5)
2 > .5*x[,1] #multiplies the first column by 0.5
3 [1] 0.5 1.0 1.5 2.0 2.5
4 > x[1,]^(1/3) #extracts cube roots from the 1st row
5 [1] 1.000000 1.259921 1.442250 1.587401 1.709976
6 > x[,1]*x[,2] #element-wise multiplies first and 2nd column
7 [1] 6 14 24 36 50
```

# Matrix Algebra

Special operators and functions exist for most of the matrix algebra. Let  $A$  and  $B$  be matrices,  $b$  a vector; then:

- ★  $A*B$ : element-wise multiplication
- ★  $t(A)$ : transpose of a matrix ( $A'$ ) or a vector ( $b'$ )
- ★  $A\%*B$ : matrix multiplication (row by column)
- ★  $A\%oB$ : outer product ( $AB'$ )
- ★  $\text{crossprod}(A,B)$ : matrix crossproduct, equivalent to  $t(A)\%*B$
- ★  $\text{diag}(A)$ : returns a vector containing the elements of the principal diagonal; it can also be used to create diagonal matrices, refer to `?diag`
- ★  $\text{solve}(A)$ : inverse ( $A^{-1}$ ) of a square matrix  $A$
- ★  $\text{solve}(A,b)$ : returns the solution of the system  $b = Ax$ , i.e.  $x = A^{-1}b$
- ★  $\text{ginv}(A)$ : Moore-Penrose Generalized Inverse of  $A$  (requires to load the MASS package)
- ★  $\text{eigen}(A)$ : returns eigenvalues and eigenvectors of  $A$

- ★ `svd(A)`: returns the Singular Value Decomposition of  $A$
- ★ `chol(A)`: returns the Choleski factorization of  $A$
- ★ `qr(A)`: returns the QR decomposition of  $A$
- ★ `cbind()`: combines matrices (vectors) horizontally, returns a matrix
- ★ `rbind()`: combines matrices (vectors) vertically, returns a matrix
- ★ `rowMeans(A)`: returns a vector of row means
- ★ `rowSums(A)`: returns a vector of row sums
- ★ `colMeans(A)`: returns a vector of column means
- ★ `colSums(A)`: returns a vector of column sums

# Working space, directories and I/O

- ★ All the objects you create are saved in a working space that can be browsed, modified and saved for future use.
  - `ls()`: plain execution of this command returns the list of the objects in the working space
  - `rm()`: this command removes an object matching its argument from the working space; say we wish to delete the object `xxx`, then we would enter `rm(xxx)` in the console. Particularly useful is the combination `rm(list=ls())`, that will remove all the objects from the working space
  - `getwd()`: execution of this command returns the current working directory. Should the need to use a different working directory arise, the `setwd()` command would change the default for the current session. Setting one's own working directory as default is particularly useful in that every input/output (I/O) procedure will then refer to it, thus making destination specifications shorter and more convenient.
  - `list.files()`: returns a character vector containing the names of the files in the current working directory
  - `save.image()`: if you want to save the current workspace to a file with extension `.Rdata`.
  - `objects.size()`: returns the (estimated) amount of memory that is

- ★ Many functions exist, that can read external files into the R working space, such as raw data (usually text files):
  - `scan()`, `read.table()` and `read.csv()`: check their help file for syntax and usage.
- ★ It is worth noting that many packages exist, among them we mention the `foreign` and the `readr` packages, that contain specific routines to import particular data formats such as, for example, those native of Excel, Minitab, S, SAS, SPSS and Stata.
- ★ For what concerns text output, R provides many ways to create `.txt` and `.csv` files:
  - `write()`, `write.table()`, `write.csv()` and `sink()`: check their help file for insights.



★ Output can also take the form of screen text and plots:

- `print()`: prints its argument on screen
- `cat()`: prints its argument by concatenation, it is useful for producing output in user-defined functions (for example, printing to screen at what step of a loop your routine is, or printing in a nice way your function results)
- `paste()`: often used as an argument to other functions, it allows to paste together, as a character string, multiple objects (useful, for example, for plot labels)
- `plot()`: the base Rfunction for plotting objects, see help

★ `jpg()`, `pdf()`, `png()`, and `postscript()`; read the related help files and remember to end every such operation with the `dev.off()` command.

# Graphical Tools

- ★ We already seen the command `plot()`. Other graphical tools include: `points()`, `lines()`, `segments()` and `abline()`, that allow you to add, respectively, points, curves, segments and straight lines to your plots
- ★ `curve()` command will allow you to plot functions of one variable over a chosen region of the  $xy$  plane
- ★ The `ggplot2` package, created by Hadley Wickham, offers a powerful graphics language for creating elegant and complex plots.  
<http://r-statistics.co/Top50-Ggplot2-Visualizations-MasterList-R-Code.html>
- ★ <http://rgraphgallery.blogspot.it/>

# Control structures, loops and miscellanea

- ★ R is in general slower than other languages, therefore a proper use of flow operators is of great importance
  - `!`, `&&`, `||`, `...`: **logical operators** allows you to evaluate the truth (logical) value of an R object, check the help for the most common ones by entering `?'`
  - **for**: for loops can be used to loop through the values of an object. The command is of the form **for** (`<index>` in `<vector>`) `<statements>`. The expressions between curly brackets are executed separately for each value in the vector

## R code

```
1 > x <- c('a','b','c')
2 > for (i in x) {
3 +   print(paste("Index=",i))
4 + }
5 [1] "Index=a"
6 [1] "Index=b"
7 [1] "Index=c"
```

- ★ For loops are very slow. You can improve code readability over more efficient solutions such as the **\*apply** functions
- ★ **if/else**: sometimes, a statement should be executed only if a certain condition is met. In this case, the **if** structure can be used, and complemented by the **else** control. It is built-up as follows: **if** (**<logical expression>**) **{<statements1>}** **else** **{<statements2>}**; **if** and **else** can be nested

## R code

```
1 > x <- 10
2 > if(x > 9) {
3 +   print("x is larger than 9")
4 + } else if(x > 7) {
5 +   print("x is larger than 7, but not larger than 9")
6 + } else {
7 +   print("x is not larger than 7")
8 + }
9 [1] "x is larger than 9"
```

- ★ `ifelse()`: If you need to set, for example, values of variables based on a vector of logical conditions, then the `ifelse` statement (which basically is a vectorized version of `if`) is the way to go. The syntax is:  
`ifelse(<condition>,{<yes statements>},{<no statements>})`

## R code

```
1 > x <- -2:2
2 > ifelse(x>=0, x, -x)
3 [1] 2 1 0 1 2
```

- ★ This is equivalent to calling `abs(x)`, i.e. the absolute value function.

★ **while**: it is one of the simplest looping structures, whose syntax is `while(<condition>) {<statements>}`. It is used to repeat <statements> while <condition> remains TRUE. For example, consider the following loop that stops as soon as a random draw from a Bernoulli distribution with parameter  $p = 0.5$  equals 1 and stores the number of draws needed to reach that point

## R code

```
1 > p <- 0.5 # probability of drawing a 1
2 > b <- 0   # result of the draw
3 > number <- 0 # number of draws
4 >
5 > while(b != 1)
6 + {
7 +   b <- rbinom(1, 1, p) # draws a Bernoulli variate
8 +   number <- number + 1 # increments the counter
9 + }
10 >
11 > print(number) #prints to screen the number of draws
```

- ★ **repeat**: this structure repeats the commands in its body until a **break** statement is reached

## R code

```
1 > k <- 0
2 > repeat {
3 +   k <- k + 1
4 +   if(k > 3) break
5 +   cat(k, "\n")
6 + }
7 1
8 2
9 3
10 >
```

- ★ `tryCatch()`: it provides a mechanism for handling with errors without breaking the loop
- ★ `system.time()`: it takes an expression as argument and returns the CPU times spent to run it
- ★ `source()`: You can save snippets of code in files other than your current one, and then source them for execution in the current session.



- ★ The R language has been originally written by statisticians
- ★ R has implemented all sort of statistical routines to perform from the simplest to the most complex task (basic or packages)

# Basic statistical functions

- ★ `summary()`: a generic function that, depending on the class of its argument outputs a summary thereof.

## R code

```
1 > x <- iris[,1]
2 > y <- iris[,2]
3 > summary(x)
4      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
5  4.300   5.100   5.800   5.843   6.400   7.900
```

## R code

```
1 > mod <- lm(y~x)
2 > summary(mod)
3 Call:
4 lm(formula = y ~ x)
5 Residuals:
6      Min       1Q   Median       3Q      Max
7 -1.1095 -0.2454 -0.0167  0.2763  1.3338
8 Coefficients:
9             Estimate Std. Error t value Pr(>|t|)
10 (Intercept)   3.41895    0.25356   13.48  <2e-16
11 x            -0.06188    0.04297   -1.44    0.152
12 ---
13 Residual standard error: 0.4343 on 148 degrees of freedom
14 Multiple R-squared:  0.01382, Adjusted R-squared: 0.007159
15 F-statistic: 2.074 on 1 and 148 DF, p-value: 0.1519
```

- ★ `mean()`: generic function for the (possibly trimmed) arithmetic mean. `colMeans()` and `rowMeans()` functions for use with matrices.
- ★ `var()`, `cov()`, `cor()`: these functions compute the variance, covariance and correlation of, say, `x` and `y` if these are vectors; `sd()` will compute the standard deviation
- ★ if the argument to `cov()` (`cor()`) is a `matrix` or a `data.frame` object, the `variance/covariance` (correlation) matrix is returned

# Probability distributions

- ★ The default R release features many classic probability distribution functions, enter `?distributions` to see them.
- ★ By convention, all functions related to a particular distribution, say `DISTR`, are encoded as follows:
  - `dDISTR()`: density function (with respect to an appropriate probability space)
  - `pDISTR()`: distribution function
  - `qDISTR()`: quantile function
  - `rDISTR()`: random variates generation

# Example: the Central Beta random variable

- ★ The related functions are then `dbeta()`, `pbeta()`, `qbeta()`, and `rbeta()`
  - `dbeta(x, shape1, shape2)`: returns the value in `x` of the density function of a Beta random variables with parameters `(shape1, shape2)`.

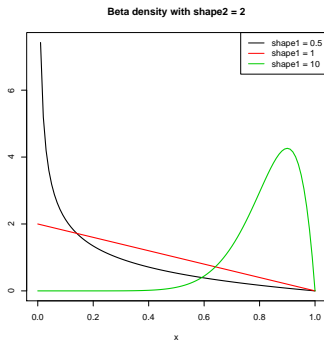
## R code

```
1 > #value of the Beta(2,2) density in x=0.5
2 > dbeta(0.5, 2, 2)
3 [1] 1.5
```

- Plot the Beta density for different combinations of its parameters:

## R code

```
1 > curve(dbeta(x,.5,2),ylab='',  
2 + main='Beta density with shape2 = 2',lwd=2)  
3 > for (i in c(1,10)) curve(dbeta(x,i,2),add=T,  
4 + col=round(i+1),lwd=2)  
5 > legend('topright',legend=c('shape1 = 0.5',  
6 + 'shape1 = 1','shape1 = 10'),col=c(1,2,11),lty=1)
```



- ★ `pbeta(q, shape1, shape2)`: returns the cumulative density function of a Beta distribution with parameters (`shape1`, `shape2`) in correspondence of quantile `q`

```
1 > beta.value <- pbeta(0.5,2,3)
2 > beta.value
3 [1] 0.6875
```

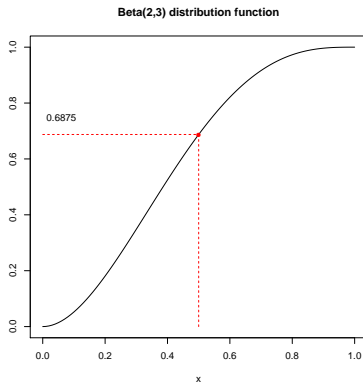
- ★ It is possible, for example, to visualize the quantile on the distribution function plot:



```

1 > curve(pbeta(x,2,3),ylab='',
2 + main='Beta(2,3) distribution function')
3 > points(0.5,beta.value,pch=16,col='red')
4 > segments(c(0,.5),c(beta.value,0),c(.5,.5),
5 + c(beta.value,beta.value),lty=2,col='red')
6 > text(.06,beta.value+.06,beta.value)

```



- ★ `qbeta(p, shape1, shape2)`: returns the quantile of order  $p$  of a Beta distribution with parameters  $(\text{shape1}, \text{shape2})$ .

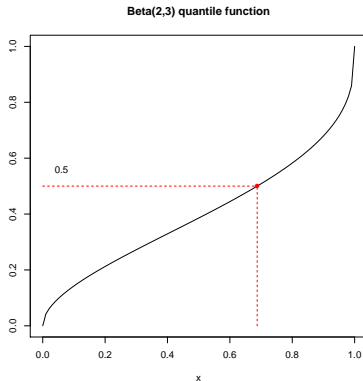
## R code

```
1 > qbeta(beta.value, 2, 3)
2 [1] 0.5
```

- ★ If we want to visualize the quantile function

## R code

```
1 > curve(qbeta(x,2,3),ylab='', main='Beta(2,3) quantile function')
2 > points(beta.value,0.5,pch=16,col='red')
3 > segments(0,0.5,beta.value,.5,beta.value,1ty=2,col='red')
4 > text(.06,0.5+.06,0.5)
```



- ★ `rbeta(n, shape1, shape2)`: generates  $n$  realizations from a Beta distribution with parameters  $(\text{shape1}, \text{shape2})$ .
- ★ How this is accomplished will be described in the next lesson. Now, imagine we wish to draw a certain number of random variates from this law and compare their distribution with the theoretical density of a Beta. First, we obtain the draws

## R code

```
1 > x <- rbeta(1000,2,3) #draws 1000 variates from Beta(2,3)
```

- ★ Check whether the sample mean and variance are close to the theoretical values. Recall that, if  $X \sim \text{Beta}(\alpha, \beta)$ , then 
$$E(X) = \frac{\alpha}{\alpha + \beta}, V(X) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)};$$
- ★ in this case, since  $X \sim \text{Beta}(2, 3)$  we would then expect  $E(X) = 0.4, V(X) \approx 0.04$

### R code

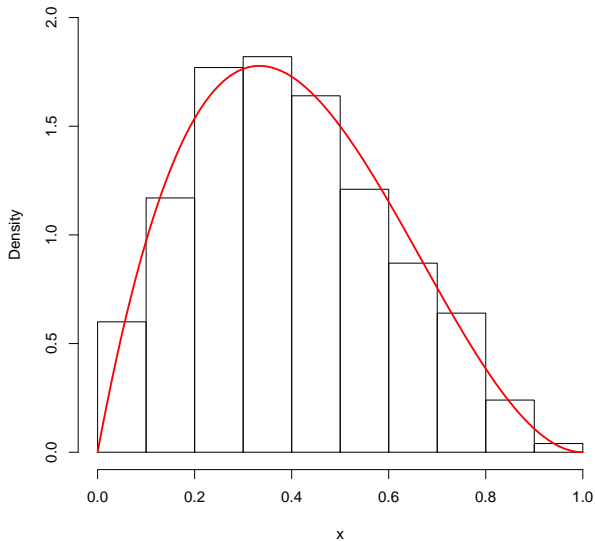
```
1 > mean(x)
2 [1] 0.4085936
3 > var(x)
4 [1] 0.04034115
```

- ★ Plot the distribution of the random variates and compare it with the true density:

## R code

```
1 > #plots the histogram
2 > hist(x,freq=F,xlim=c(0,1),ylim=c(0,2),
3 + main='Random draws from a Beta(2,3)\n and true density')
4 > #overimposes the theoretical density
5 > curve(dbeta(x,2,3),
6 + col='red',add=T,lwd=2)
```

### Random draws from a Beta(2,3) and true density



- ★ Several functions exist, that can handle the most common statistical models; e.g.:
  - `lm()`: fits classic linear models, and provides a wealth of statistics on residuals, as well as hypothesis testing on the model's parameters
  - `glm()`: fits generalised linear models, and takes as input the description of the linear predictor and of the error term; again, the most common model-related statistics are computed as a by-product of the routine
  - `anova()`: based on one or more previously fitted model objects, computes the related analysis of variance (or deviance) tables.
- ★ There are also routines for survival analysis, longitudinal models, time series, multilevel modeling, partial least squares, structural equations, ...!



# Hypothesis testing

- ★ The base R release contains functions for many of the most common statistical tests, such as, for example:
  1. `t.test()` (for one- and two-sample  $t$ -test),
  2. `ks.test()` (for the Kolmogorov-Smirnov one- and two-sample tests),
  3. and `wilcox.test()` (for the Wilcoxon-Mann-Whitney tests)
- ★ Some testing routing will be used during Lessons 2 and 3.

# Statistical plots

- ★ `plot()`: produces graphical representations based on the class of its argument
- ★ `hist()`: it plots an histogram based on a vector of observations
- ★ `boxplot()`: produces box-and-whisker plot(s) of the given (grouped) values, based on a `data.frame` or `list` object
- ★ `density()`: computes kernel density estimates based on a vector of observations
- ★ `acf()`: computes and plots estimates of the autocovariance or autocorrelation function based on a vector of observations. Useful when working with time series, stochastic processes, and when evaluating a random number generator

- ★ `pairs()`: plots the pairwise scatterplots based on a `matrix` or `data.frame` object, useful for a first visual inspection of a dataset, or to summarize dependence between pairs of variables
- ★ `image()`: creates a grid of rectangles with colors corresponding to the values in a given `matrix` object (3-dimensional data)
- ★ `contour()`: creates a contour plot, or add contour lines to an existing plot

# Writing R functions

- ★ One of the most appealing features of R, is that everybody can create new functions and packages to be made available to all other R users;
- ★ Most of R is actually written in R (apart from a few routines that are written in C for efficiency reasons), which means that you could basically take every function and suitably modify it to meet your needs
- ★ An R function is defined by an assignment such as:

## R code

```
1 MyFunction <- function(arg1, arg2, ...) {  
2   expression  
3   ...  
4   expression  
5   value  
6 }
```

- ★ **expression** is an R statement (possibly involving other functions) that uses some of the arguments **arg1, arg2, ...** to compute a value, that is the outcome of the function
- ★ The curly braces indicate the beginning and the end of the function; it is also possible to set default values for the arguments
- ★ a function should always **output a value**, since everything that happens within those curly braces (unless otherwise specified) is **local** to the function itself and, as such, temporary

# Example 1

★ Create a function that:

1. takes two vectors as arguments
2. checks if they are of the same length, otherwise returns an error message

## R code

```
1 length.check <- function(x,y) {#x and y are arg1 and arg2
2
3   if (length(x)==length(y)){#the logical condition to be met
4
5     print('Same length.')#value condition is TRUE
6
7   } else {
8
9     print('Different length.')#value if condition is FALSE
10
11   }
12 }
```

## ★ Test our function

### R code

```
1 > a <- 1:10
2 > b <- letters[2:11]
3 > c <- runif(15)
4 > length.check(a,b)
5 [1] "Same length."
6 > length.check(a,c)
7 [1] "Different length."
8 > length.check(b,c)
9 [1] "Different length."
```

## Example 2

★ We now wish to write a function that:

1. has two arguments, of which the second one defaults to the value 1
2. computes the  $k$ -th sample moment of the first object (suppose it's a vector of sample observations), where  $k$  is defined by the second argument
3. outputs the computed value

### R code

```
1 #By setting y=1 in the arguments ,
2 #we define its default value
3 k.moment <- function(x, y=1) {
4   value <- mean(x^y)      #stores the result
5
6   value      #outputs the result
7
8 }
```



## ★ Test our function

### R code

```
1 > a <- 1:10
2 > k.moment(a,1)
3 [1] 5.5
4 > k.moment(a)
5 [1] 5.5
6 > k.moment(a,2)
7 [1] 38.5
8 > k.moment(a,3)
9 [1] 302.5
```

## Example 3: probability law not already implemented in R

### ★ Reciprocal Distribution

- continuous random variable with bounded support  $[a, b]$  ( $0 < a < b$ ) and density  $f(x; a, b) = [x \ln(b/a)]^{-1}$

### ★ Code for:

- density function
- distribution function: it is easily obtained by integration and equals  $F(q; a, b) = \ln(q/a)[\ln(b/a)]^{-1}$
- quantile function: it is  $Q = F^{-1}$

## R code

```
1  #density function
2
3  drec <- function(x,a,b) {
4
5    #checks constraints
6
7    stopifnot(x>=a&& x<=b&& 0<a&&a<b)
8
9    #outputs density value
10
11    (x*log(b/a))^(-1)
12
13 }
```

## R code

```
1 #distribution function
2
3 prec <- function(q,a,b) {
4
5   #checks constraints
6
7   stopifnot(0<a&&a<b)
8
9   #if q<=a, returns 0, otherwise
10  #if q>=b, returns 1, otherwise
11  #returns cdf value
12
13  ifelse(q<=a,0,
14        ifelse(q>=b,1,
15              log(q/a)/log(b/a)))
16
17 }
```

## R code

```
1 #quantile function
2 qrec <- function(p,a,b) {
3
4     #checks constraints
5
6     stopifnot(p>=0&& p<=1&& 0<a&&a<b)
7
8     #computes q numerically by
9     #finding the zero of 'F(q)-p'
10    #in [a,b], and returns it
11
12    uniroot(
13      function(x) prec(x,a,b)-p,
14      lower=a, upper=b)$root
15
16 }
```

## R code

```
1 > a <- c(1,2.5,5)
2 > drec(a,1,10)
3 [1] 0.4342945 0.1737178 0.0868589
```

## R code

```
1 > a <- seq(0,1,.3)
2 > qrec(a,10,15)
3 [1] 10
4 Warning messages:
5 1: In if (is.na(f.lower)) stop("f.lower = f(lower) is NA") :
6 the condition has length > 1
7 and only the first element will be used
8 2: In if (is.na(f.upper)) stop("f.upper = f(upper) is NA") :
9 the condition has length > 1
10 and only the first element will be used
11 3: In if (f.lower * f.upper > 0)
12 stop("f() values at end points not of opposite sign") :
13 the condition has length > 1
14 and only the first element will be used
```

- ★ `qrec()` internally calls `uniroot()`, that does not accept vector-valued functions as argument
- ★ to solve this, it suffices to `Vectorize` the function

#### R code

```
1 > qrec <- Vectorize(qrec)
2 > qrec(a,10,15)
3 [1] 10.00000 11.29347 12.75426 14.40397
```



# Exercises

- ★ Using the `seq()` function, construct the following sequences:
  1. integer values from 1 to 10
  2. 32 equi-spaced values from 1 to 10
  3. values distant 0.15 from each other, ranging from 1 to 10; how many of them are there? (Tip: use `length()`).
- ★ Using the `replicate()` function, obtain a vector of 1000 realizations of means from uniform samples of size 100 each. Plot the histogram of the standardized observation and superimpose a Standard Normal curve (tip: set `freq=FALSE` when plotting the histogram).
- ★ Read the help for the `iris` dataset and describe the variables it contains. Create a two-pages .pdf document file with:
  1. the pairwise scatter plots for all the involved variables, with different colours indicating different species on the first page
  2. the plot of the kernel density estimate of `Petal.Length`, regardless of the species and the histogram of `Sepal.Length` for the `versicolor` species only on the second page (tip: `?density`).

- ★ Read the help for the `iris` dataset and describe the variables it contains. Create a two-pages .pdf document file with:
  1. the pairwise scatter plots for all the involved variables, with different colours indicating different species on the first page
  2. the plot of the kernel density estimate of `Petal.Length`, regardless of the species and the histogram of `Sepal.Length` for the `versicolor` species only on the second page (tip: `?density`).
- ★ Write a function that:
  1. accepts two arguments `a` and `b`, that must be non-negative integer with default value 1
  2. plots the function  $f(x; a, b) = |\cos(e^{-ax^2+bx})|$  on the range  $x \in [-1, 1]$
  3. returns the value of  $\int_{-1}^1 f(x; a, b)dx$  on screen and saves it in a .txt file in the current working directory.

