

1 R Basics

Exercise 1 (Data structures and subsetting)

- (a) Create a `list` named `list1` with the following components

`x` : a numeric vector 35, 34, 33, ..., 7, 6, 5, 6, 7, ..., 33, 34, 35 (don't type the single numbers!)

`y` : a factor with the first five letters of the Latin alphabet as levels, each repeated 12 times, except the last one which should be repeated 13 times

`mat` : a numeric 4×4 matrix whose entries are randomly drawn from an exponential distribution with $\lambda = 5$. Use the `set.seed()` function to make your result reproducible.

`list2` : a list with the components

- `t` : a numeric variable with the value 35
- `d` : a data.frame with the variables
 - `gender` : a factor with elements male, male, male, female, female, female, male, male, male, female, female, female
 - `age` : a numeric vector with elements 23, 48, 37, 37, 19, 54, 21, 20, 41, 26, 35, 32

- (b) Use subsetting operators to access the following information

- (a) the 4th and the 7th element of the vector `x`
- (b) the entry on position (3, 4) of the matrix `mat`
- (c) the first six elements of the first column of the data frame `d`
- (d) the age of the female individuals

What is the difference between the usage of `[]`, `[[]]` and `$` ?

- (c) Modify the created objects in the following ways

- (a) redefine the levels of the factor `y` from lowest to highest as 'c', 'd', 'b', 'a', 'e'
- (b) eliminate the first row and the third column of `mat`
- (c) add a column `age2` to `d` with a factor taking the value 'old', if the individual's age is above the threshold `t`, and the value 'young' otherwise
- (d) exclude the individuals that are at least 50 and strictly younger than 21 from the study

Exercise 2 (Loops and data manipulation)

- (a) Create a matrix X of dimension 2500×12 containing random numbers $x_{ij} \sim \mathcal{N}(\mu = 0, \sigma^2 = 3)$.
- (b) Create a new matrix $(y_{ij})_{ij}$ based on the one defined above, in which each column is standardized, i.e. $y_{ij} = \frac{x_{ij} - \bar{x}_j}{\sigma_j}$, where \bar{x}_j is the sample mean of the j -th column, and σ_j its sample standard deviation.

- (c) From X keep only the rows for which at least 6 out of the 12 values are greater than 0. Hint: Remember that a logical vector can be converted to a numeric vector where `FALSE` takes the value 0, and `TRUE` the value 1.
- (d) Check, if the elements of the 3rd column of `mat` from Exercise 1 are within the range defined by the elements in the first two columns. The output should be a logical vector.
- (e) In each row of X replace the maximum value with the minimum value.
- (f) Write a for-loop to do the following calculation for $i = 1, \dots, 2500$:
- if the i -th element of the first column X is positive, then add to the subsequent element $x_{i+1,j}$ a random number $u \sim \mathcal{U}(-1, 1)$
 - otherwise, if x_{ij} is negative, add to the subsequent element $x_{i+1,j}$ a random number $u \sim \mathcal{U}(-2, 2)$.
- (g) Consider the data.frame resulting from

```
set.seed(2015)
w <- runif(10)
x <- runif(10)
y <- runif(10)
z <- runif(10)
dat <- data.frame(a = w, b = x, c = y, d = z).
> dat
```

	a	b	c	d
1	0.06111892	0.70285557	0.6919100	0.75185674
2	0.83915986	0.39172125	0.4067718	0.25684487
3	0.29861322	0.03306277	0.2109400	0.38967137
4	0.03143242	0.40940319	0.6652073	0.88448520
5	0.13857171	0.74234713	0.7377556	0.57390551
6	0.35318471	0.88301877	0.9190050	0.18367673
7	0.49995552	0.26623321	0.8734601	0.11168811
8	0.07707116	0.07427093	0.8012774	0.32047459
9	0.65134483	0.81368426	0.5243978	0.09095567
10	0.51172371	0.38194719	0.1272213	0.47959896

Now imagine in each row the entries define two intervals $[a, b]$ and $[c, d]$. Add a column to `DF` with entry `TRUE`, if the intervals overlap and `FALSE`, if they don't.

Find ways to avoid for-loops in (a)–(c). Useful functions might be

`rnorm()`, `runif()`, `apply()`, `mean()`, `sd()`, `scale()`, `rowSums()`, `min()`, `max()`.

If you are not familiar with some of these functions, use `?function` to check the documentation. Compare the `system.time()` of the solutions with and without for-loops.

Exercise 3 (Basic graphical tools)

- (a) Consider again the matrix X from Exercise 2 (a) and create the following plots:
- A scaled histogram of the first column. Also add a dashed red line representing the estimated density.
 - A quantile-quantile plot comparing the sample quantiles of the first and second column. Also add a line with intercept 0 and slope 1 to improve comparability. Try different `pch`-values and colours.
 - Boxplots of the first, fourth and seventh row in one plotting window. Rename the group labels to `blue`, `green` and `yellow`.

- (b) Load the data set `MathAchieve` from the package `nlme`. In two neighboring plotting windows create boxplots of the achievements (`MathAch`) for the groups `Male` and `Female` (left plot) and a scatterplot of the achievement in dependence of the `SES` (socio-economic status) together with the according linear regression line in red (right plot).

2 Advanced Graphics

Exercise 1 (From basic plots to complex graphics)

Read the data set `school_math.raw`. It describes the mathematical achievements (`mathach`) of male and female (`Sex`) students at 4 different schools (`school`). Additionally there is information about the socio-economic status (`ses`) of the students family and the sector (`Sector`) of the school (public or catholic).

- (a) Get familiar with the data set, e.g. via the functions `summary()` or `head()`.
- (b) Create the following four plots in only one plotting window:
- a **scatterplot** of the achievement vs. the socio-economic status. Add a line representing the estimated linear relationship between those variables.

a **histogram** of `mathach`.

a **scaled histogram** of `ses`. Add a line with the estimated density.

boxplots of the mathematical achievement for each school separately.

First, don't modify any options. In a second step include the following variations:

- The color of the numbers on the axes should be red.
- Data points should be represented by filled squares (■) instead of circles (○).
- Use squared plotting regions for each plot.
- All occurring lines should be dashed.

Modify your plot by changing the settings within the single plot functions and in the overall graphics options using `par()`.

Hint: Make a backup of the default `par()` settings, such that you can reset the settings to default afterwards.

Hint: The documentations of `plot`, `points`, `lines`, `abline`, `par`, `hist`, `boxplot` might be helpful.

- (c) Reset the graphic settings to the default values.

Exercise 2 (The layout function)

Add additional information to the scatterplot from Exercise 1 in form of boxplots of the data at the axes. For this purpose, get familiar with the function `layout()` and its options. Also use `layout.show()` to visualize different settings.

Exercise 3 (Lattice plots and grouped data)

- (a) Estimate regression lines of the form

$$\text{mathach} = \beta_0 + \beta_1 \cdot \text{ses} + \epsilon$$

for the different schools separately. Plot your point estimates against the school. Use one plotting window for the intercepts and another one for the slopes. The y -axes should be labeled β_0 and β_1 , respectively.

Hint: Look at the function `lmList()` from the package `nlme`. Use `expression()` for the labeling.

- (b) Plot the mathematical achievement (in dependence of the `ses`) of boys as green triangles and of girls as red squares. Add a legend to the plot indicating this grouping.
- (c) Load the packages `lattice` and `nlme`. For the school data set, we want to compare the data in the different groups. For that purpose create a `groupedData` object with response `mathach`, covariate `ses` and `school` as grouping variable. Now
 - Plot the new object. What is the difference between the results from `plot(data)` and `plot(data_new)`, where `data` is the original data set without grouping structure and `data_new` is the `groupedData` object?
 - Create box-whiskers plots (`bwplot()`) of the residuals from an overall linear model according to `mathach ~ ses` grouped by the different schools.
 - Introduce a random intercept for each school and estimate a mixed model using the function `lme()`. Compare the results with your findings from (a). Try `compareFits()` and `comparePred()`.
 - Plot the confidence intervals of the estimated intercepts and slopes for the different schools.

Exercise 4 (3D images)

Load the package `plot3D` and the data set `sambia.raw`. Visualize the relationships between `zscore`¹ and `breastfeeding`² and `age`³ both individually and jointly using the function `scatter3D()`.

3 Data Management

Exercise 1 (Spreadsheet-like data)

- (a) Load the file `knee.txt` using the command `read.table()`. Compare what happens when using `header = TRUE` and `header = FALSE`.
- (b) What class do the single variables have in R? Are the assigned classes reasonable in all cases? Define the classes of `th`, `gen` and `pain` as factor directly when reading in the data.
- (c) Additionally rename the variables to *treatment*, *age*, *sex* and *agony*.
- (d) Open the data set `knee2.txt` first in a text editor. What do you observe? Load the data into R such that missing values are automatically identified as those.
- (e) Do the same with `knee3.txt`. What is the difference in contrast to `knee2.txt`?
- (f) Add a column to the data frame `knee2` with the value *old*, if `age` is larger than 40 and *young* otherwise. Save the resulting data set in the following formats:
 - `*.Rdata`,
 - `*.raw` without quotes and with `:` as a separator,
 - `*.dat` without column or row names,
 - `*.csv`.

For the latter compare the results from `write.table`, `write.csv` and `write.csv2`.

¹Here, the `zscore` measures the nutrition status of children in Zambia. Highly negative values indicate severe under nutrition.

²Duration of breastfeeding in month

³Age of the child in month

Exercise 2 (Handling date objects and ordering data in R)

The dataset `movies.csv` contains information on the 10 most successful movies of each for the years 2012-2015.

- (a) What is the class of the variables `release` and `end`? Obviously those two variables are the dates of the release and the last day of the movies in cinemas. Use `as.Date()` to transform them into dates.
- (b) Add a column to the data frame with the number of days the movies have been shown in cinemas. Add another column with the number of complete weeks.
- (c) Add a column with the year of release.
- (d) Rank the movies according to the number of weeks the movies have been shown in cinemas.
- (e) (optional) Write a piece of code which has as output the name of each year's most successful movie (according to its gross income).
- (f) (optional) Find out on which `weekday()` Carl Friedrich Gauss was born.

Exercise 3 (Converting objects and *.RData-objects)

Load the data stored in `data.Rdata`. Use `ls()` to find out which objects are stored in that file. Convert the list into a vector and the data frame into a matrix. What do you observe? Try to explain what happens.

Exercise 4 (Foreign data structures)

- (a) The file `blutdruck.sav` is a SPSS data file. Use the according function from the package `foreign` to load it into R. What do you need to do to directly create a data frame in R?
- (b) Use the function `read.sas7bdat()` from the package `sas7bdat` to read in the SAS file `s05_01.sas7bdat`.
- (c) Read in the Stata file `golf.dta` with `read.dta()`.

Exercise 5 (Reshaping data, optional exercise)

Consider the dataset `airquality`, which is already stored in R. It is a dataset in wide format, i.e. each variable has its own column.

- (a) Load the package `reshape2` and convert the dataset into long format. Don't modify any options yet.
- (b) For this dataset it might be interesting to have an overview about the values of the different climate variables `Ozone`, `Solar.R`, `Wind` and `Temp` for each day of the month. Use the option `id.vars` to order the data according to `Month` and `Day`. Also change the name of the variable column to `climate_var` while reshaping.
- (c) Often, the situation is the other way around. You find a dataset looking like the one you created in (2). Now you are interested in comparing (or find relationships between) the different climate variables directly. Hence, you want your dataset look somehow like the original `airquality` data. Use the function `dcast()` to split the column `climate_var` into single columns for each variable.

4 Functions, Debugging & Condition Handling

Exercise 1 (Simple functions and restrictions)

In this exercise, we want to build a function to calculate the binomial coefficient of two integers n, k

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} \quad (1)$$

Here, we will treat the binomial coefficient as not defined for negative numbers and non-integers⁴. Follow the next steps carefully:

- (a) What are the definition and basic properties of the factorial of an integer? Also recall the properties of (1) for two integers.
- (b) Write a simple function `my_factorial()` which calculates the factorial of an integer n . Don't add any restrictions on n yet. Also, evaluate `my_factorial` at numbers from $\mathbb{R} \setminus \mathbb{N}$. What do you observe?
- (c) Let your function return an error message, if n is not an integer. The built-in function `is.integer()` will not help you in this case (why?). Hence, write your own function which returns `TRUE` for integers and `FALSE` for non-integers. Check if your function fulfills the properties from (a). If not, add necessary exceptions.
- (d) Write a function `my_binomial()` to calculate the expression in (1). Make use of your results from (b) and (c). Evaluate the expression `my_binomial(2, 3)`. Is the resulting error message reasonable? Use `traceback()` to locate where the error occurred.
- (e) Finally, add an exception for the case $k > n$ to `my_binomial()`. The function should then return the value 0. All other restrictions should be kept as before.

Exercise 2 (Scope and environments)

For a given vector $x = (x_1, \dots, x_n)$, the function

```
running_mean = function(k) {  
  1 / k * sum(x[1:k])  
}
```

calculates the running mean of order $k \leq n$, i.e.

$$\bar{x}_k = \frac{1}{k} \sum_{j=1}^k x_j$$

Guess the output of the following two scenarios. Then, check if your prediction was correct.

- (a)

```
x = 1:5  
running_mean = function(k) {  
  1 / k * sum(x[1:k])  
}  
running_mean(2)
```
- (b)

```
x = 1:5  
rm(list = ls())  
running_mean = function(k) {  
  1 / k * sum(x[1:k])  
}  
running_mean(2)
```

⁴Note that there are extensions based on the Γ -function for these cases.

- (c) Explain what went wrong in (b). Modify the function `running_mean` to make it independent from the environment/workspace.
- (d) Add reasonable restrictions on x and k .

Exercise 3 (Closures and environments)

In contrast to problems caused by different environments as in the previous exercise, one can explicitly use these environments.

1. Write a function `n.root` that allows for generating the family of n -th root functions

$$\sqrt[n]{x} = x^{\frac{1}{n}}$$

For that purpose write a function within a function, such that x is used in the inner environment and n in the outer environment.

2. Can you use `n.root` to explicitly calculate the roots of a number?
3. Based on `n.root` write functions which calculate the square and cubic root, respectively.
4. Use `n.root` to calculate the first 10 roots of $x = 500$.

Exercise 4 (Default values and optional arguments, optional exercise)

For a sample $x = (x_1, \dots, x_n)$ and some $\epsilon \geq 0$, the Windsorized mean is defined as follows. If \bar{x} and s are the sample mean and standard deviation, respectively, replace the values in the sample which are larger than $\bar{x} + \epsilon s$ by $\bar{x} + \epsilon s$ and the values in the sample which are smaller than $\bar{x} - \epsilon s$ by $\bar{x} - \epsilon s$. The Windsorized mean is then given by

$$\mu_W^x = \bar{y},$$

where y is the transformed sample.

- (a) Write a function that computes the Windsorized mean of a sample. Include an argument `eps`, where y is the transformed sample. which corresponds to the ϵ as explained above with 2 as default value.
- (b) Make your function more general and allow for calculating the running Windsorized mean⁵ and use the result from Exercise 2. Find a way to modify the function `running_mean`, such that the default is the calculation of the classical mean of the full sample.

5 Object Oriented Programming

Exercise 1 (S3-class objects and methods)

- (a) Write a function called `BestTeam()` with the following output:
 - a named `list` with the name of your favourite sports team and a logical object `football` indicating whether it is a football team or not;
 - the resulting object should be of S3 class `FavoriteTeam`.

Hint: For the latter, use the function `append()`, such that the class is a vector consisting of the original class (`list`) and the newly defined one. If necessary have a look at `?class`.

- (b) Define a method `setFavoriteTeam` which adapts the favorite team of a `FavoriteTeam` class object. It should trigger an error, if it is used with an object of any other class.
- (c) What is the output of

⁵i.e. the running mean of the Windsorized sample and not the Windsorized mean of the truncated sample

```
myTeam = BestTeam(isFootball = FALSE, team = "EWE Baskets Oldenburg")
summary(myTeam)
```

and why?

Write an extension to the generic method `summary()` for `FavoriteTeam` class objects, which produces something like

```
sport = other
team = EWE Baskets Oldenburg
```

The possible values of `sport` should be `other` and `football`, depending on the logical indicator defined in (a).

Exercise 2 (S4-class objects and methods)

Write a function that creates an S4 class `FavoriteTeam2`. The resulting object should be similar to the one in Exercise 1. Also add an extension to the generic `summary` function and a method `setFavoriteTeam`.

Hint: Helpful functions are `setClass`, `setGeneric` and `setMethod`.

Exercise 3 (Object specific advanced graphics)

Reconsider Exercise 2 from the tutorial on advanced graphics (use of the `layout` function). Create an appropriate S4 class `newData`. Also write an extension to the function `plot`, such that for each object `x` of the new class the command

```
plot(x)
```

produces the before mentioned plot.

6 Profiling, Performance & Parallelization

Exercise 1 (Loops and if statements) Recall Exercise 2 (c) from the first section.

(a) Use `system.time()` or `proc.time()` to compare the solutions resulting from

- (i) `if ... else` statements within a `for` loop,
- (ii) `ifelse()` within a `for` loop,
- (iii) when calculating the row sums within a `for` loop,
- (iv) `rowSums` as presented in the solution of exercise sheet 1.

Check whether all solutions lead to the same results (see `?identical`). What are your conclusions (if any)?

(b) Install and load the package `microbenchmarking` and read the documentation of the likewise named function.

- (i) Compare again the four approaches from (a).
- (ii) Load the package `ggplot2` and use `autoplot()` to visualize the findings.

Exercise 2 (R specific tricks - implemented functions)

(a) Consider the least squares estimator

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

from the linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

Load the data from `linear.RData` and compare the computation times of

- using `%*%` and `solve` and
 - using `crossprod` whenever possible and inverting $\mathbf{X}^\top \mathbf{X}$ via a QR decomposition.
- (b) Use the function `microbenchmarking` again to compare the computation times of `mean(y)` and `sum(y)/length(y)`.

Exercise 3 (Parallelization)

Install and load the platform specific packages for parallelization (e.g. `parallel` and `doParallel` for Windows) as well as the data set `iris`.

1. Exclude the species `setosa` from the data set and use `glm` to estimate the influence of the `Sepal.Length` on (the probability of observing one of) the `Species` in a logit model.
2. Parallelize (use `detectCores()`, `registerDoParallel()`, `foreach()`, `%dopar%` for Windows) a bootstrap procedure to approximate the distributions of the estimators (intercept and covariate effect).
3. Compare the computational time for $B = 1000$ and $B = 10000$ bootstrap samples, respectively, with the time needed without parallelization.

7 R Interaction with Other Languages

Exercise 1 (Getting started with Rcpp)

- (a) Run the functions `meanC` and `sumC` of the lecture, to check whether `Rcpp` is working on your computer.
- (b) Translate the below defined R-functions in C++. These functions create the elements of the Fibonacci sequence either as recursive functions or as vectors. Benchmark your results.

```
fibR_recursive <- function(n) {  
  if (n == 0) return(0)  
  if (n == 1) return(1)  
  return (fibR_recursive(n - 1) + fibR_recursive(n - 2))  
}  
  
fibR_vector <- function(n) {  
  if (n == 0) fibo <- 0  
  if (n == 1) fibo <- c(0,1)  
  if (n > 1) {  
    fibo <- rep(0,times=n+1)  
    fibo[2] <- 1  
    for(i in 3:(n+1)) {  
      fibo[i] <- (fibo[i-1] + fibo[i-2])  
    }  
  }  
  return(fibo)  
}
```

- (c) Implement from scratch a function to estimate the variance of a vector in R and C++. Compare the computational time of both functions.

- (d) Implement a function in R and C++, which sorts a vector in increasing order. Compare the computational time of both functions. An example of a sorting algorithm is the insertion sorting algorithm as defined below in pseudo code. Here $A=(A[1], \dots, A[n])$ is the vector to be sorted.

```

INSERTIONSORT(A)
  for i <- 2 to length(A) do
    value_to_be_sorted <- A[i]
    j <- i
    while j > 1 and A[j-1] > value_to_be_sorted do
      A[j] <- A[j - 1]
      j <- j - 1
    A[j] <- value_to_be_sorted

```

8 Numerics and Simulation

Exercise 1 (Cancellation)

Use `curve()` to visualize the function

$$f(x) = \frac{1 - \cos(x)}{x^2}$$

- (a) between -15 and 15 as well as
- (b) between -4×10^{-8} and 4×10^{-8} .

Try to explain what you observe.

Exercise 2 (A Monte-Carlo approximation of π) Consider a circle of radius 1 around $(0, 0)$. The circle consists of the points $(x, y) \in \mathbb{R}$ with the restriction

$$x^2 + y^2 \leq 1 \tag{2}$$

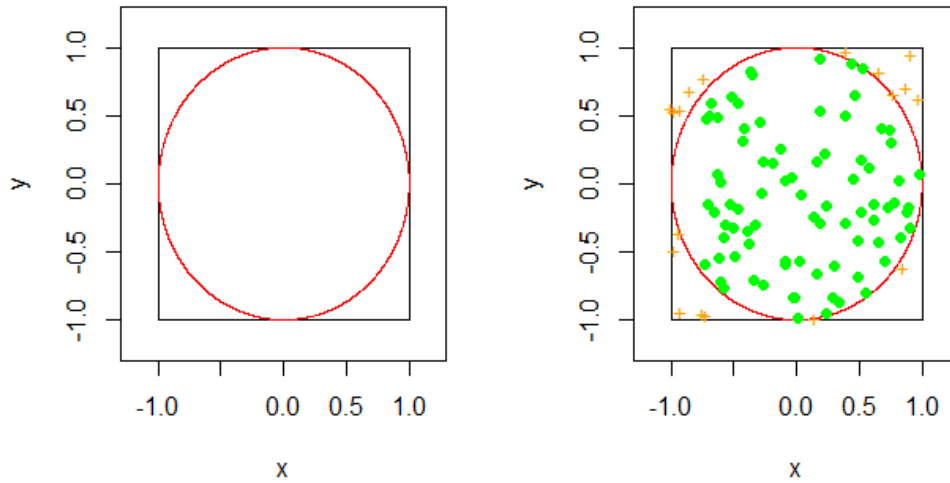
such that its area is exactly π . We use this property to find a suitable approximation of π . Translate the following ideas into R-Code:

- Embed the mentioned circle into a square with vertices $(-1, -1)$, $(1, -1)$, $(1, 1)$ and $(-1, 1)$ (see Figure 1), the ratio of the areas of the circle and the square is

$$R = \frac{\pi}{4} \tag{3}$$

- Approximate R to find an approximation of π .
- Proceed as follows:
 1. Draw random points from within the square (use `runif()`).
 2. Identify those points that lie in the circle (Equation (1) might be helpful).
 3. Use the proportion of the points from 2. as approximation for R .
 4. Use Equation (2) to approximate π .
 5. Start with $n = 10$ random points and increase the sample size as long as the relative error of the approximation is larger than 10^{-5} (make use of `while`).

Figure 1: *Left:* A circle within a square. *Right:* Random points.



How large is the required sample size in your case? Compare your results with those from other course participants.

Exercise 3 (Propagation of round off errors - the Vancouver stock exchange bug)

The Vancouver stock market was established in January 1982 with a starting index of 1000. After 22 month of trade, it decreased to a value of about 525, whereas the expected value was above 1000. What happened?

- After each transaction the index was recalculated
- The value of each transaction was available up to four decimal places
- For the calculation of the new index, only three were used (the fourth was cut off and not rounded)

Approximate the overall error under the following assumptions:

- 2900 transaction per day
- 20 business days per month
- 22 month of trade
- an average error of 0.00045 per transaction

Based on the approximation, what is the true index after those 22 month of trade?

Exercise 4 (Simulation and optimization) The density of a Gaussian distributed random variable X with mean μ and variance σ^2 is

$$f_{\mu, \sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

Hence for a sample x_1, \dots, x_n from the distribution of X , the log-likelihood with respect to μ (up to additive constants) is given by

$$l(\mu) = -\frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \mu)^2.$$

- (a) Write a function to compute the log-likelihood of μ for a given sample. Consider the variance σ^2 to be known.
- (b) Simulate one sample of size $n = 500000$ with $\mu = 3$, $\sigma = 1$ and use `optimx::optimx` to compare the performance of different optimization routines when calculating the Maximum-Likelihood estimator.
- (c) Generate a list of 100 samples of size $n = 200$ and compute the Maximum-Likelihood estimator for each of the samples with an optimization routine of your choice. Make use of `lapply`.
- (d) Visualize the results.

9 Building R Packages

Exercise 1 (Create your own package) Create an R package based on at least one of the following files (which are available via GitHub)

1. `beta.hat.R`: cf. tutorial 6 (Profiling); include the data stored in `linear.RData`
2. `my binomial.R`: cf. tutorial 4 (Functions)
3. `NicePlot.R`: cf. tutorials 3 (Graphics) and 5 (Object Orientation); include the data stored in `NicePlot.RData`
4. `pi approx.R`: cf. tutorial 8 (Numerics and Simulations)

In all cases, write short documentations for the occurring functions and data sets. For including the latter, `devtools::use_data()` might be helpful. Also check and build your package. You are free to work in groups.