

## R FUNCTIONS FOR ECONOMISTS

A table of data in R is called »data frame«. Such a table can be arbitrarily manipulated and analyzed with R functions and R commands. An example of such a table, which we call mydata, is given below:

	ID	Age	Height	Gender	Place
1	1	30	180	1	LJ
2	2	25	170	0	LJ
3	3	50	175	1	MB
4	4	40	170	0	LJ
5	5	16	190	1	KP
6	6	22	178	1	NM
7	7	29	177	0	GO

Table is manipulated by `mydata[ , ]`. The notation before the comma refers to rows, the notation after the comma refers to columns. Examples:

```
mydata[c(-4, -5), ]
mydata[, c(1, 2, 3)] or mydata[, c(1:3)]
mydata[4, 3] <- 172
```

In the first example, we delete the 4th and 5th rows from the data table, in the second example, we select the 1st, 2nd and 3rd variables (columns), and in the third example, we replace the element in the 4th row and the 3rd column (i.e. 170) with the value 172. The single variable is accessed with the `$` character. Example: `mydata$Age`.

The data table can be ordered according to the selected variable with the function `order`:

```
mydata[order(mydata$Age), ]
mydata[order(-mydata$Height), ]
```

In the first case, the data table is arranged in ascending order according to the age of the persons, in the second case in descending order according to the height of the persons.

The character `<-` (assign left) assigns a function or the command on the right to the selected object on the left.

The `%>%` (pipe) sign means »then«. To use it, we need to install the library `dplyr`.

R functions and commands are written in R Markdown, and all text must be enclosed in quotes. The R program distinguishes between upper and lower case, punctuation, brackets, etc., but it is not sensitive to spaces.

Wide format of the data table

	ID	Trust_A	Trust_B
1	1	8	3
2	2	5	5
3	3	4	6
4	4	7	2
5	5	2	9

Long format of the data table

	ID	Trust	Factor
1	1	8	A
2	2	5	A
3	3	4	A
4	4	7	A
5	5	2	A
6	1	3	B
7	2	5	B
8	3	6	B
9	4	2	B
10	5	9	B

Data specified in a frequency table should be converted to raw data using the `rep` function (see function explanation below):

Frequency table

	Members	Frequency
1	1	2
2	2	2
3	3	3
4	4	1

Raw data

	Members	Frequency
1	1	2
1.1	1	2
2	2	2
2.1	2	2
3	3	3
3.1	3	3
3.2	3	3
4	4	1

BASIC R FUNCTIONS and COMMANDS	
<b>install.packages ()</b> <b>library ()</b>	Installing and activating the new library.
<pre>install.packages("psych") library(psych)</pre> <p>Basic functions for installing and activating new libraries that are not yet installed in the R program. The new library needs to be installed only once, after that it should only be activated. Basic libraries (e.g. base, stats) are installed and activated automatically.</p>	
<b>help ()</b>	Help in using functions and commands.
<pre>help(read.table)</pre> <p>It explains how each function or command works and how it should be used. Help can also be accessed by clicking F1 when the name of the function is selected with the mouse in R Markdown.</p>	
<b>read.table ()</b>	Reading a .csv data file.
<pre>mydata &lt;- read.table("./data.csv",                     header = TRUE,                     sep = ";",                     dec = ",")</pre> <p>Reading a .csv file that contains variable names in the first line. We specify how the variables are separated (sep) and the type of decimal separator (dec).</p>	
<b>read_xlsx ()</b>	Reading an Excel file with data.
<pre>library(readxl) mydata &lt;- read_xlsx("./data.xlsx") mydata &lt;- as.data.frame(mydata)</pre> <p>Reading files saved in Excel format. To use the function, we need to install and activate the readxl library. When we import data, we also use the as.data.frame function to convert the object (table with data) into data.frame format.</p>	
<b>data.frame ()</b>	Manual data entry creating a data table.
<pre>mydata &lt;- data.frame("Age" = c(30, 25, 50),                     "Height" = c(180, 170, 175))</pre> <p>We have created a data table with three units and two variables.</p>	

<b>c ()</b>	Creating a vector of elements.
<pre>c(4, 5, 10) c("M", "M", "F")</pre> <p>If you specify more than one element, write c before the parenthesis. Letters/texts must always be enclosed in quotation marks.</p>	
<b>sort ()</b>	Sorting values by size.
<pre>sort(mydata\$Age, decreasing = FALSE)</pre> <p>We sort the values of the variable in ascending or descending order.</p>	
<b>factor ()</b>	Creating factor.
<pre>factor(mydata\$Gender,        levels = c(1, 0),        labels = c("M", "F"))</pre> <p>Categorical variables, whose values are usually encoded with numeric codes, are converted into a factor. Only then will they be treated correctly as categorical variables by the R program. In the argument levels we write down the codes with which the variable is entered in the data table, and in the argument labels we write down what these codes mean.</p> <pre>mydata\$GenderF &lt;- factor(mydata\$Gender,                         levels = c(1, 0),                         labels = c("M", "F"))</pre> <p>We have created a new variable GenderF and stored it in an existing data table.</p>	
<b>head ()</b> <b>print ()</b> <b>tail ()</b>	Display a table of data.
<pre>head(mydata) print(mydata) tail(mydata)</pre> <p>We display the first 6 rows of the data table (head), the entire table (print) or the last 6 rows of the data table (tail).</p>	

## R FUNCTIONS FOR ECONOMISTS

<p><b>round()</b> <b>floor()</b> <b>ceiling()</b></p> <pre>round(mydata, 2)</pre> <p>The values in the data table are rounded to two decimal places.</p> <pre>floor(214.8)</pre> <p>Round the value down to the first integer. The result of the function in this particular case is 214.</p> <pre>ceiling(214.8)</pre> <p>Round the value up to the first integer. The result of the function in this particular case is 215.</p>	<p>Rounding values.</p>	<p><b>scale()</b></p>	<p>Standardization</p> <pre>mydata\$Height_z &lt;- scale(mydata\$Height)</pre> <p>Standardization of the individual variables. In the concrete case, we assign the new variable to the existing data table.</p> <pre>mydata_std &lt;- scale(mydata)</pre> <p>Data table standardization.</p>
<p><b>ifelse()</b></p> <pre>ifelse(mydata\$Age &gt; 22,        yes = 1,        no = 0)</pre> <p>A function is used to execute a conditional command. In the concrete example, the result of the function is a vector of values, where people older than 22 years have the value 1, otherwise 0.</p> <pre>ifelse(mydata\$Age &gt;= 22,        yes = 1,        no = 0)</pre> <p>A function is used to execute a conditional command. In the concrete example, the result of the function is a vector of values, where people older or equal to 22 years have the value 1, otherwise 0.</p>	<p>Conditional command.</p>	<p><b>sqrt()</b></p>	<p>Square root.</p> <pre>sqrt(mydata\$Height)</pre> <p>Calculation of root values.</p>
<p><b>seq()</b></p> <pre>seq(from = 0, to = 100, by = 20)</pre> <p>The function is used to create a vector of values with a constant difference between them. In the concrete case, the result is <code>c(0, 20, 40, 60, 80, 100)</code>.</p>	<p>A sequence of values.</p>	<p><b>aggregate()</b></p>	<p>Aggregation.</p> <pre>aggregate(mydata\$Height, by = list(mydata\$Gender),           FUN = sum)</pre> <p>Aggregation of a variable with respect to a categorical variable. In this particular case, we add up all the heights of the individuals for each gender separately. Enter any function (sum, mean etc.) in the argument FUN to calculate the parameter of interest.</p>
<p><b>log()</b> <b>exp()</b></p> <pre>mydata\$lnHeight &lt;- log(mydata\$Height) mydata\$lnHeight &lt;- log(podatki\$Height, base = 10)</pre> <p>Calculation of the logarithm. In the first case we calculate the natural logarithm, in the second case the decimal logarithm.</p>	<p>Calculation of logarithm and exponential function (anti-logarithmization).</p>	<p><b>recode()</b></p>	<p>Value recoding.</p> <pre>library(car) recode(mydata\$Height,        "150:160 = 155; 160:170 = 165; 170:180 = 175; 180:190 = 185")</pre> <p>Recoding a range of values to a new value. In this particular case, all values between 150 and 160 cm are recoded to 155 cm, and so on. (Threshold values are recoded to the value which they first fall according to the written code; in this particular case, the value 160 is recoded to the value 155, not 165. If we were to write the code in the order "160:170 = 165; 150:160 = 155...", the value 160 would be recoded to the value 165). To use the function, we need to install and activate the library <code>car</code>.</p>
		<p><b>rep()</b></p>	<p>Reproduction of value.</p> <pre>mydata[rep(1:nrow(mydata), times = mydata\$Frequency), ]</pre> <p>With the function we reproduce each row in the data table as many times as the value of frequency (argument <code>times</code>). The function is used, for example, when we have data in sorted form (frequency table – aggregated data) and we need to convert it to unsorted form – raw data.</p>

## R FUNCTIONS FOR ECONOMISTS

<b>replace_with_na()</b>	Replacing the value with NA.
<pre>library(dplyr) library(naniar) mydata &lt;- mydata %&gt;%   replace_with_na(replace = list(parliament = c(77, 88, 99),                                 police = c(77, 88, 99),                                 politics = c(77, 88, 99)))</pre>	
Replace numeric codes representing missing values with the value NA. To use the function, we need to install and activate the library <code>dplyr</code> and <code>naniar</code> .	
<b>drop_na()</b>	Removing units with missing values.
<pre>library(tidyr) mydata &lt;- drop_na(mydata)</pre>	
Removing units that have a missing value for any of the variables. To use the function, we need to install and activate the library <code>tidyr</code> .	
<b>order()</b>	Sorting units by values.
<pre>print(mydata[order(-mydata\$Points, mydata\$Time), ])</pre>	
We display a table of data sorted in descending order by the number of points scored, and units with the same number of points scored are sorted in ascending order by writing time	

## DESCRIPTIVE STATISTICS

<b>mean() median() sd() var() range() min() max() sum()</b>	Estimation of statistical parameters: mean, median, standard deviation, variance, range, minimum, maximum, sum.
<pre>mean(mydata\$Age) median(mydata)</pre>	
Evaluation of a single parameter for a selected variable or the entire table.	
<b>quantile()</b>	Estimation of quantile.
<pre>quantile(mydata\$Age, p = 0.35)</pre>	
Using the argument <code>p</code> , which determines the rank of the quantile, we can estimate any quantile. In this particular case, we estimate the 35th centile.	
<b>summary()</b>	Estimation of parameters.
<pre>summary(mydata)</pre>	
The function provides an estimate of the minimum, first quartile, mean, median, third quartile, and maximum. In the case of variables defined as a factor, the function gives the frequency of each category.	
<b>stat.desc()</b>	Estimation of parameters.
<pre>library(pastecs) stat.desc(mydata) round(stat.desc(mydata), 1)</pre>	
The function provides a wide range of parameter estimates. The values are given to 7 decimal places, so it is useful to combine the function with the function <code>round</code> . To use the function, we need to install and activate the <code>pastecs</code> library.	
<b>describe()</b>	Estimation of parameters.
<pre>library(psych) describe(mydata)</pre>	
The function gives a wide range of parameter estimates. To use the function, we need to install and activate the <code>psych</code> library.	

## R FUNCTIONS FOR ECONOMISTS

<b>describeBy ()</b>	Estimation of parameters, separated by groups.
<pre>library(psych) describeBy(mydata\$Age, group = mydata\$Gender)</pre> <p>The function provides a wide range of parameter estimates separated by categories (groups). To use the function, we need to install and activate the <code>psych</code> library.</p>	
<b>sapply ()</b>	Estimation of selected parameter.
<pre>sapply(mydata, FUN = mean) sapply(mydata, FUN = quantile, probs = c(0.01, 0.99))</pre> <p>The function gives estimates of the selected parameter for all variables located in the data table.</p>	
<b>table ()</b>	The number of repetitions of single value.
<pre>table(mydata\$Gender)</pre> <p>The function specifies the number of repetitions of a single value for the selected variable.</p>	
<b>nrow ()</b> <b>length ()</b>	Number of rows. The length of the vector.
<pre>nrow(mydata) length(mydata\$Age)</pre> <p>The number of rows in the data table or length of the vector (variable).</p>	
<b>freq ()</b>	Frequency table.
<pre>library(sjmisc) table(mydata\$Height)</pre> <p>Frequency table for the selected numerical variable. To use the function, we need to install and activate the <code>sjmisc</code> library.</p>	

## STATISTICAL DISTRIBUTIONS

<b>pnorm ()</b>	Calculation of the area under the normal distribution curve.
<pre>pnorm(q = 5, mean = 3, sd = 2, lower.tail = TRUE)</pre> <p>The function gives the area under the probability density curve for the normal distribution to the left of the selected value of the variable (<code>lower.tail = TRUE</code>) or to the right of the selected value (<code>lower.tail = FALSE</code>).</p>	
<b>qnorm ()</b>	Calculate the value of the variable for the selected area for the normal distribution.
<pre>qnorm(p = 0.30, mean = 3, sd = 2, lower.tail = TRUE)</pre> <p>The function gives the value of the variable for the selected area under the probability density curve for a normal distribution. The selected area can be located to the left (<code>lower.tail = TRUE</code>) or to the right of the calculated limit value (<code>lower.tail = FALSE</code>).</p>	
<b>rnorm ()</b>	Generating random numbers from a given normal distribution.
<pre>rnorm(n = 20, mean = 3, sd = 2)</pre> <p>We generate a selected number (n) of random values from a given normal distribution.</p>	
<b>dnorm ()</b>	Calculation of the density for the normal distribution at the selected value of the variable.
<pre>dnorm(x = 2.5, mean = 3, sd = 2)</pre> <p>We calculate density function at the selected value of the normal distribution.</p>	
<b>qt ()</b>	Calculation of the value of the variable for the selected area for the t-distribution.
<pre>qt(p = 0.05, df = n-1)</pre> <p>The function gives the value of the variable in the selected area under the density curve for the t distribution at the selected degrees of freedom (df).</p>	

## SAMPLING

### choose ()

Number of combinations without repetition.

```
library(combinat)
choose(6, 4)
```

From a population of size 6, we select all possible samples of 4 units without repetition. In this particular case, it is  $\binom{6}{4} = 15$ . To use the function, we need to install and activate the library `combinat`.

### combn ()

Print of all possible samples.

```
library(combinat)
combn(mydata, 4)
```

Print the values of the variable for each sample. To use the function, we need to install and activate the library `combinat`.

```
combn(mydata, 4, mean)
```

Output of mean estimates for all possible samples.

## VIZUALIZATION OF DATA

### hist ()

Histogram.

```
hist(mydata$Age,
     main = "Distribution of age",
     ylab = "Frequency",
     xlab = "Age",
     breaks = seq(from = 20, to = 50, by = 5))
```

We display the data with a histogram. The main argument gives the title of the graph, `ylab` and `xlab` give the axis names, and the width of the columns is given by the `breaks` argument.

### boxplot ()

Boxplot.

```
boxplot(mydata$Age)
```

We plot the data using a boxplot. It consists of a rectangle defining the first, second, and third quartiles and a vertical line defines the minimum and maximum of the variables under study. Potential outliers are marked with circles.

### ggplot ()

Tool for graphical representation of data.

```
library(ggplot2)
ggplot(mydata, aes(x = Age)) +
  geom_histogram(binwidth = 5, colour = "gray") +
  ylab("Frequency")
```

```
library(ggplot2)
ggplot(mydata, aes(y = Age)) +
  geom_boxplot()
```

```
library(ggplot2)
ggplot(mydata, aes(y = Height, x = Age)) +
  geom_point()
```

We select the table with the data, and with the argument `aes` we specify which variable we draw on each axis. Then we build the chart using `+` signs. The argument `geom_` specifies the chart type. To use the function, we need to install and activate the `ggplot2` library. A more detailed explanation can be found in the document `ggplot2.pdf`.

### scatterplot ()

Scatterplot.

```
library(car)
scatterplot(mydata$Height ~ mydata$Age,
            smooth = FALSE,
            boxplots = FALSE,
            regLine = FALSE,
            ylab = "Height in cm", xlab = "Age in years")
```

Scatter plot for a pair of numerical variables. The first variable is on the y-axis, the second on the x-axis. With the arguments `smooth = FALSE` we turn off smoothing, with the argument `boxplots = FALSE` we turn off the display of the boxplot on the axes, with the argument `regLine = FALSE` we turn off the linear regression function, and with the arguments `ylab` and `xlab` we name the axes. To use the function, we need to install and activate the library `car`.

```
library(car)
scatterplot(mydata$Height ~ mydata$Age | mydata$Gender,
            smooth = FALSE,
            boxplots = FALSE,
            ylab = "Height in cm", xlab = "Age in years")
```

Scatter plot, separated by categorical variable.

## R FUNCTIONS FOR ECONOMISTS

### scatterplotMatrix()

Matrix of scatterplots.

```
library(car)
scatterplotMatrix(mydata, smooth = FALSE)
Matrix of scatter plots for each pair of numeric variables. Non-numeric variables must be excluded from the analysis. To use the function, we need to install and activate the library car.
```

## HYPOTHESIS TESTING

### t.test()

Hypothesis about the arithmetic mean.  
Hypothesis about the difference between two arithmetic means.

```
t.test(mydata$Age,
       mu = 25,
       alternative = "two.sided")
```

Testing of the hypothesis about the value of the population arithmetic mean. The value of the null hypothesis is entered into the `mu` argument, and the `alternative` argument is used to determine whether it is a two sided test ("two.sided") or one sided test: ("less") or ("greater").

```
t.test(mydata$Weight1, mydata$Weight2,
       paired = TRUE,
       alternative = "less")
```

Testing of the hypothesis about the difference between two arithmetic means for the dependent samples (`paired = TRUE`).

```
t.test(mydata$Height ~ mydata$Gender,
       paired = FALSE,
       var.equal = FALSE,
       alternative = "greater")
```

Testing of the hypothesis about the difference between two arithmetic means for independent samples (`paired = FALSE`). With the argument `var.equal = FALSE` we choose the Welch correction for the possible difference of the group variances.

### anova\_test()

Analysis of variance for the dependent samples, rANOVA.

```
library(rstatix)
anova_test(dv = Trust,
           wid = ID,
           within = System,
           data = mydata_long)
```

We define the dependent variable (`dv`), we define the unit identifier (`wid`), we define the factor indicating what the dependent variable measures (`within`), and the data must be in a long format. To use the function, we need to install and activate the `rstatix` library.

### aov()

Analysis of variance for the independent samples, ANOVA.

```
aov(Salary ~ Country,
     data = mydata)
```

We define the dependent variable and the factor that separates the units into groups.

### prop.test()

Hypothesis about the proportion.  
Hypothesis about the equality of two proportions.

```
prop.test(x = 150, n = 250,
         p = 0.50,
         correct = FALSE,
         alternative = "two.sided")
```

We determine the number of events (`x`) and the number of all trials (`n`). Using the argument `p`, we determine the null hypothesis, exclude the correction, and determine the alternative hypothesis.

```
prop.test(c(195, 45), c(234, 89),
         correct = TRUE,
         alternative = "two.sided")
```

Testing of the hypothesis about the equality of two proportions. In the specific case, we check whether the proportion 195/234 is statistically different from the proportion 45/89.

### shapiro.test()

Hypothesis about whether the variable is normally distributed.

```
shapiro.test(mydata$Difference)
```

We use the Shapiro-Wilk test to check whether the variable is normally distributed.

 FUNCTIONS FOR ECONOMISTS

<p><b>pairwise.t.test()</b></p> <pre>pairwise.t.test(mydata\$Salary, g = mydata\$Country,                 paired = FALSE,                 p.adj = "bonf")</pre> <p>Testing the hypothesis about the difference of two arithmetic means for dependent samples (paired = TRUE) or independent samples (paired = FALSE) for each pair of variables. A Bonferroni correction is often used when calculating <math>p</math>-values.</p>	<p>Comparison of all pairs of arithmetic means with <math>t</math>-tests.</p>
<p><b>cohens_d()</b> <b>interpret_cohens_d()</b></p> <pre>library(effectsize) cohens_d(mydata\$Age, mu = 25) interpret_cohens_d(0.52, rules = "sawilowsky2009")</pre> <p>We determine the effect size and then interpret it based on various rules. To use the function, we need to install and activate the effectsize library.</p>	<p>Effect size, determined using Cohen's D statistic.</p>
<p><b>wilcox_test()</b></p> <pre>wilcox_test(mydata\$Age,             mu = 23,             correct = FALSE)</pre> <p>The Wilcoxon test for median (nonparametric test) is calculated by determining the median (<math>\mu</math>) and excluding the correction.</p> <pre>wilcox_text(mydata\$Weight1, mydata\$Weight2,             paired = TRUE,             correct = FALSE, exact = FALSE,             alternative = "two.sided")</pre> <p>The Wilcoxon test for equality of two distribution locations for dependent samples (Wilcoxon signed rank test - nonparametric test) is computed by specifying both variables and choosing the argument paired = TRUE. We specify the alternative hypothesis and we exclude the correction and the calculation of the exact <math>p</math>-value.</p> <pre>wilcox_text(mydata\$Height ~ mydata\$Gender,             paired = FALSE,             correct = FALSE, exact = FALSE,             alternative = "less")</pre>	<p>Wilcoxon test for median.</p> <p>Wilcoxon test for equality of locations of two distributions for dependent and independent samples.</p>

<p><b>binom.test()</b></p> <pre>binom.test(x = 150, n = 250,            p = 0.70,            alternative = "two.sided")</pre> <p>We determine the number of events (<math>x</math>) and the number of all trials (<math>n</math>). Using the argument <math>p</math>, we determine the null hypothesis and specify the alternative hypothesis. If <math>p = 0.50</math>, we perform a sign test.</p>	<p>Binomial test (proportion test). Sign test.</p>
<p><b>friedman_test()</b></p> <pre>library(rstatix) friedman_test(Trust ~ System   ID,               data = mydata_long)</pre> <p>We define a dependent variable and a factor that indicates what the dependent variable measures, and then an identifier that assigns the measurement to an individual unit. The data must be in long format. This is a nonparametric test. To use the function, we need to install and activate the rstatix library.</p>	<p>Friedman ANOVA.</p>
<p><b>kruskal.test()</b></p> <pre>kruskal.test(Salary ~ Country,              data = mydata)</pre> <p>We define the dependent variable and the factor that separates the units into groups. This is a nonparametric test.</p>	<p>Kruskal-Wallis Rank Sum test.</p>



## R FUNCTIONS FOR ECONOMISTS

<p><b>summarySE ()</b></p> <pre>library(Rmisc) summarySE(mydata,           measurevar = "Height",           groupvars = "Gender",           conf.interval = 0.95)</pre> <p>We define data containing a numeric variable for which we calculate a confidence interval for the arithmetic mean (<code>measurevar</code>), and a categorical variable that divides the values of the variable into groups (<code>groupvars</code>). To use the function, we need to install and activate the <code>Rmisc</code> library.</p>	<p>Calculation of the parameters of the confidence intervals for the arithmetic mean by groups.</p>
<p><b>VarTest ()</b></p> <pre>library(DescTools) VarTest(mydata\$Height,         sigma.squared = 5,         alternative = "two.sided")</pre> <p>Testing the hypothesis about the value of the variance. With the argument <code>sigma.squared</code> we determine the value in the null hypothesis, and with the <code>alternative</code> argument we determine whether it is a two-sided hypothesis (<code>"two.sided"</code>), a hypothesis that is directed to the left (<code>"less"</code>) or to the right (<code>"greater"</code>). To use the function, we need to install and activate the <code>DescTools</code> library.</p> <pre>library(DescTools) VarTest(mydata\$Height ~ mydata\$Gender,         alternative = "two.sided")</pre> <p>Hypothesis about the assumption of equality of two variances (whether their ratio is equal to 1). In the specific case, we check whether the height variances differ between the sexes, and the <code>alternative</code> argument is used to determine the alternative hypothesis. To use the function, we need to install and activate the <code>DescTools</code> library.</p>	<p>Hypothesis about the value of the variance. Hypothesis about the ratio of two variances.</p>

## ANALYSIS OF CATEGORICAL VARIABLES

<p><b>chisq.test ()</b></p> <pre>chisq.test(mydata\$Gender, mydata\$Place,            correct = FALSE)</pre> <p>We use the test to check if there is a relationship between two categorical variables. If both variables have exactly two categories, we use the argument <code>correct = TRUE</code>.</p> <pre>chisq.test(x = c(150, 200, 350), p = c(0.30, 0.30., 0.40),            correct = FALSE)</pre> <p>We test whether two distributions are the same. In the argument <code>p</code> we write the expected probabilities, which are used to determine the expected frequencies.</p>	<p>Pearson <math>\chi^2</math>-test. <math>\chi^2</math>-test for given probabilities.</p>
<p><b>fisher.test ()</b></p> <pre>fisher.test(mydata\$Gender, mydata\$Place)</pre> <p>We use the test to check if there is a relationship between two categorical variables.</p>	<p>Fisher's exact probability test.</p>

## CORRELATION AND REGRESSION ANALYSIS

**cor()**  
**cor.test()** | Estimation of the correlation coefficient.  
Hypothesis about the value of the correlation coefficient.

```
cor(mydata$Age, mydata$Height,
    method = c("pearson", "spearman"))
```

Estimation of the correlation coefficient between two numerical variables. Either Pearson's or Spearman's correlation coefficients can be calculated.

```
cor(mydata)
```

Estimation of correlation matrix. The data table must contain only numerical variables.

```
cor.test(mydata$Age, mydata$Height,
    method = c("pearson", "spearman"))
```

Hypothesis about the value of the correlation coefficient.

**rcorr()** | Correlation matrix

```
library(Hmisc)
rcorr(as.matrix(mydata), type = c("pearson", "spearman"))
```

An estimate of the correlation matrix, containing either Pearson's or Spearman's correlation coefficients for each pair of numeric variables.  $p$ -values are also provided below the matrix to test the hypothesis about the value of the correlation coefficient. To use this function, the library `Hmisc` must be installed and activated.

**ppcor()** | Estimation of the partial correlation coefficient.

```
library(ppcor)
ppcor(mydata)
```

An estimate of the partial correlation coefficient for each pair of numeric variables in the data table. To use this function, the library `ppcor` must be installed and activated.

**lm()** | Linear regression model estimation (OLS method).

```
lm(Height ~ Age + Gender + Age:Gender,
    data = mydata)
```

Estimation of the linear regression function using the least squares method. The dependent variable is indicated to the left of the tilde (~), while the explanatory variables are indicated to the right, separated by the sign +. Interactions are included by indicating the sign between the variables :

```
fit <- lm(Height ~ Age + Gender,
    data = mydata)
summary(fit)
```

Showing the results of the estimated regression function.

**glm()** | Linear regression model estimation (ML method).

```
glm(Height ~ Age + Gender,
    data = mydata)
```

Estimation of the linear regression function using the maximum likelihood method. The dependent variable is indicated to the left of the tilde (~), while the explanatory variables are indicated to the right, separated by the sign +. Interactions are included by indicating the sign between the variables :

```
glm(Smoker ~ Age + Gender,
    family = binomial,
    data = mydata)
```

Estimation of a binary logistic regression function.

```
fit <- glm(Smoker ~ Age + Gender,
    family = binomial,
    data = mydata)
summary(fit)
```

Showing the results of the estimated regression function.

**ols\_test\_breusch\_pagan()** | Breusch-Pagan heteroskedasticity test.

```
library(olsrr)
ols_test_breusch_pagan(fit)
```

Test for the presence of heteroskedasticity. To use the function, we need to install and activate the `olsrr` library.

## R FUNCTIONS FOR ECONOMISTS

<b>shapiro.test()</b>	Hypothesis if the variable is normally distributed.
<pre>shapiro.test(mydata\$StdResiduals)</pre> <p>We use the Shapiro-Wilk test to check whether the variable is normally distributed.</p>	
<b>vif()</b>	Checking for multicollinearity.
<pre>vif(fit)</pre> <p>Checking the degree of correlation between explanatory variables.</p>	
<b>lm.beta()</b>	Estimation of standardized partial regression coefficients.
<pre>library(lm.beta) lm.beta(fit)</pre> <p>Estimation of the standardized partial regression coefficients of the regression model. To use the function, we need to install and activate the library <code>lm.beta</code>.</p>	
<b>anova()</b>	Comparison of two regression models in terms of fit to data.
<pre>anova(fit1, fit2)</pre> <p>Statistical comparison of two regression models estimated by the least squares method.</p> <pre>anova(fit1, fit2, test = "Chi")</pre> <p>Statistical comparison of two regression models estimated by the maximum likelihood method.</p>	
<b>lm_robust()</b>	Linear regression model estimation with robust standard errors.
<pre>library(estimatr) lm_robust(Height ~ Age + Gender,           se_type = "HC1",           data = mydata)</pre> <p>Estimation of the regression model with robust White's standard errors. The correction is applied when the assumption of homoskedasticity is violated. To use this feature, the library must be installed and activated <code>estimatr</code>.</p>	

<b>lme()</b>	Estimation of a multilevel linear regression model (ML method).
<pre>library(nlme) fit &lt;- lme(Stress ~ Age + GenderF + SizeF,           random = ~ 1   ID_hospital / ID_department,           method = "ML",           data = mydata)</pre> <p>Estimation of multilevel (hierarchical) regression model using the maximum likelihood method. We specify the regression model, and the argument <code>random</code> is used to determine the random regression constant. In a specific case, this is determined at the hospital and department level. To use the function, we need to install and activate the <code>nlme</code> library.</p>	

## PRINCIPAL COMPONENT ANALYSIS AND FACTOR ANALYSIS

<b>cortest.bartlett()</b>	Bartlett's test of sphericity.
<pre>library(psych) cortest.bartlett(R, n = nrow(mydata))</pre> <p>Bartlett's test of sphericity of the correlation matrix. The input element is the value of the correlation matrix created with the <code>cor()</code> function. We also specify the number of units in the sample. To use the function, the <code>psych</code> library must be installed and activated.</p>	
<b>det()</b>	Determinant of the correlation matrix.
<pre>det(R)</pre> <p>Determinant of the correlation matrix, with which we check the degree of correlation between the variables.</p>	
<b>KMO()</b>	Kaiser-Meyer-Olkin measure of sampling adequacy.
<pre>library(psych) KMO(R)</pre> <p>Calculation of KMO statistics and individual MSA statistics to check the adequacy of each variable. To use this function, the library <code>psych</code> must be installed and activated.</p>	

## R FUNCTIONS FOR ECONOMISTS

<p><b>PCA ()</b></p> <pre>library(FactoMineR) library(factoextra) PCA(mydata,      scale.unit = TRUE,      ncp = 3,      graph = FALSE)</pre> <p>Applying the principal component analysis to selected numeric variables. With the argument <code>scale.unit = TRUE</code> we perform the normalization of the variables, with the argument <code>ncp</code> we determine the number of components. To use the function, we need to install and activate the libraries <code>FactoMineR</code> and <code>factoextra</code>.</p> <pre>pca &lt;- PCA(mydata,            scale.unit = TRUE,            graph = FALSE)</pre> <p><code>get_eigenvalue(pca)</code></p> <p>We show the eigenvalues of the principal components.</p> <pre>fviz_eig(pca,          choice = "eigenvalue",          addlabels = TRUE)</pre> <p>Eigenvalue chart.</p> <pre>fviz_pca_var(pca, repel = TRUE)</pre> <p>Graphic representation of the loadings of the principal components.</p> <pre>fviz_pca_biplot(pca)</pre> <p>Graphical representation of the position of the units in relation to the values of the first two principal components.</p>	<p>Principal Component Analysis.</p>
<p><b>fa.parallel ()</b></p> <pre>library(psych) fa.parallel(mydata,            sim = FALSE,            fa = c("pc", "fa"))</pre> <p>Parallel analysis to determine the number of principal components or factors to be retained in the final result. We include only numeric variables in the data table, exclude simulation, and use the argument <code>fa</code> to indicate whether we are performing principal component (<code>pc</code>) or factor (<code>fa</code>) analysis. To use the function, we need to install and activate the <code>psych</code> library.</p>	<p>Parallel analysis.</p>

<p><b>corPlot ()</b></p> <pre>library(psych) corPlot(R)</pre> <p>Graphical representation of the correlation matrix. The input element is the value of the correlation matrix created by the <code>cor ()</code> function. The colors indicate the strength of the correlation for each pair of variables. To use the function, we need to install and activate the <code>psych</code> library.</p>	<p>Graphical representation of the correlation matrix.</p>
<p><b>fa ()</b></p> <pre>library(psych) library(GPARotation) fa(mydata,    covar = FALSE,    nfactors = 3,    fm = "minres",    rotate = "oblimin",    inpute = "mean")</pre> <p>Performing factor analysis for selected numerical variables. The argument <code>covar</code> is used to specify whether the analysis is performed on the covariance matrix (<code>covar = TRUE</code>) or on the correlation matrix (<code>covar = FALSE</code>), the argument <code>nfactors</code> is used to specify the number of factors, the argument <code>fm</code> is used to specify the factor method (recommended use of <code>minres</code> – minimum residual method), the argument <code>rotate</code> is used to specify the type of factorial rotation (oblique <code>oblimin</code> or orthogonal <code>varimax</code>). Missing values can be replaced by average values of variables (<code>inpute = "mean"</code>). To use this function, the libraries <code>psych</code> and <code>GPARotation</code> must be installed and activated.</p> <pre>factor &lt;- fa(mydata,             covar = FALSE,             nfactors = 3,             fm = "minres",             rotate = "oblimin",             inpute = "mean")</pre> <pre>print.psych(factor,             cut = 0.3,             sort = TRUE)</pre> <p>We present the results of the factor analysis. Low factor loadings are not shown (<code>cut</code>), and factor loadings are shown in descending order (<code>sort</code>).</p> <pre>fa.diagram(factor)</pre> <p>Graphic representation of the factor model.</p>	<p>Factor Analysis.</p>

## R FUNCTIONS FOR ECONOMISTS

### alpha ()

Cronbach alpha.

```
library(psych)
alpha(mydata, check.keys = TRUE)
```

Calculation of Cronbach's alpha for selected numeric variables. The argument `check.keys = TRUE` ensures that all indicators point in the same direction. To use this function, the `psych` library must be installed and activated.

## CLUSTERING

### get\_dist ()

Calculation of distances between units.

```
library(factoextra)
get_dist(mydata, method = "euclidean")
```

Calculation of Euclidean distances between units. The distance type is specified with the argument `method`. The distances "manhattan" and "minkowski" are also commonly used. To use this function, the `factoextra` library must be installed and activated.

```
distances <- get_dist(mydata, method = "euclidean")
fviz_dist(distances)
```

Graphical representation of the distance matrix.

### get\_clust\_tendency ()

Calculation of Hopkins statistics.

```
library(factoextra)
get_clust_tendency(mydata,
                   n = nrow(mydata)-1,
                   graph = FALSE)
```

Calculation of Hopkins statistics, which is used to check whether the data are suitable for carrying out the clustering.

### hclust ()

Hierarchical clustering.

```
library(dplyr)
library(factoextra)
WARD <- mydata %>%
  get_dist(method = "euclidean") %>%
  hclust(method = "ward.D2")
```

Implementation of hierarchical clustering based on selected numerical variables. In the first step, the data is selected, then (`%>%`) the distances between the units are calculated (`get_dist`), and then hierarchical clustering (`hclust`) is performed based on these distances according to the selected algorithm (`method`). It is recommended to use Ward's algorithm, which is combined with squared Euclidean distance (`ward.D2`). To use the function, the libraries `dplyr` and `factoextra` must be installed and activated.

```
fviz_dend(WARD)
```

Representation of dendrogram (classification tree).

```
cutree(WARD, k = 3)
```

Clustering into a selected number of groups.

### hkmeans ()

K-means clustering.

```
library(factoextra)
hkmeans(mydata,
        k = 3,
        hc.metric = "euclidean",
        hc.method = "ward.D2")
```

Clustering by the k-means method, starting from the solution obtained by hierarchical clustering. We need to specify the number of groups (`k`) into which we want to group the units, and the arguments `hc.metric` and `hc.method` should match the selected algorithm for hierarchical sorting, based on which we decided on the number of groups. To use the function, we need to install and activate the `factoextra` library.

```
kmeans <- hkmeans(mydata,
                  k = 3,
                  hc.metric = "euclidean",
                  hc.method = "ward.D2")
```

```
fviz_cluster(kmeans)
```

Graphical representation of the clustering, with the first and second principal components shown on the axes.

## TIME SERIES

### lm()

Calculation of linear or exponential trend (OLS method).

```
lm(Overnights ~ t,
   data = mydata)
```

Linear trend function. The variable `t` is the time, defined as 1, 2, 3, ..., N.

```
lm(lnOvernights ~ t,
   data = mydata)
```

Exponential trend function. The variable `t` is the time, defined as 1, 2, 3, ..., N.

### ts()

Determination of time series.

```
ts(mydata$Overnights,
   start = c(2020, 1),
   end = c(2022, 12),
   frequency = 12)
```

We define the variable as a time series. The `start` argument is used to determine the first observation (in the concrete case it is 2020, the first month), the `end` argument is used to determine the last observation (in the concrete case it is 2022, the last month), and the `frequency` argument is used to determine how many time units are within a period. In the concrete case, the data is given monthly, so a total of 36 observations.

```
overnights <- ts(mydata$Overnights,
                 start = c(2020, 1),
                 end = c(2022, 12),
                 frequency = 12)
```

```
plot(overnights,
     ylab = "Monthly number of overnights")
```

Graphical display of time series.

### decompose()

Decomposition of time series.

```
decompose(Overnights,
          type = "multiplicative")
```

Breakdown of the time series into basic elements (trend, cyclical component, period, and irregular component). The variable we analyze must be the result of the `ts()` function. For the type of decomposition, we choose the multiplicative approach.

```
overnights <- ts(podatki$Overnights,
                 start = c(2020, 1),
                 end = c(2022, 12),
                 frequency = 12)
```

```
decomposition <- decompose(overnights,
                           type = "multiplicative")
```

```
plot(decomposition)
```

Graphical representation of a time series decomposition.