

EFR summary

Programming (AAC), FEM71003
2025-2026



Lectures 1 to 7
Weeks 1 to 7

Deloitte.

DeNederlandseBank
EUROSYSTEEM

Details

Subject: Programming 2025–2026

Teacher: P Breuer

Date of publication: 13.12.2024

© This summary is intellectual property of the Economic Faculty association Rotterdam (EFR). All rights reserved. The content of this summary is not in any way a substitute for the lectures or any other study material. We cannot be held liable for any missing or wrong information. Erasmus School of Economics is not involved nor affiliated with the publication of this summary. For questions or comments contact summaries@efr.nl

Koffietje doen?

Start jouw carrière bij BDO

Maak kennis met BDO Accountants & Adviseurs,
de beste plek om als toptalent aan de slag te gaan.
De koffie staat voor je klaar. Vertellen wij je over
wat jij kan bijdragen, en jij ons over je ambities.



*Scan de QR-code
en plan jouw
koffiemoment in.*



werkenbijbdo.nl ►

BDO

Programming – lecture 1

Learning objectives

You will learn how to:

- Write R scripts to import and explore data
- Use data tidying techniques and visualization tools
- Apply data transformation techniques
- Perform regression analysis and write basic loops and functions
- Apply string manipulation techniques and regular expressions
- Employ textual analysis techniques

How to leverage AI for learning programming:

- Code explanations
- Debugging assistance
- Practice problems
- NOT replace your own coding skills, AI is here to assist.

R and RStudio: coding basics

Learning to program is like learning a new language, it takes practice and patience. Programming skills are transferable to other analytics tools (tools like Java, SQL, Python etc.)

What is **R**?

- Programming language for statistical computing and graphics.
- Open source and free software

What is **RStudio**?

- Integrated development environment (IDE) for R.
- Open source and free software
- Helpful features
 - o User-friendly interface, code autocompletion and exhaustive help on any object.

In the interface of RStudio there are 4 windows that are aligned on the screen as follows:

1. Script editor Enter commands and save to file	3. Workspace View active objects and history of commands executed
2. Console Enter commands and see output	4. Files/Plots/Packages/Help Open files, view plots, install and load packages or use the help function.

You can use R as calculator, by simply typing the calculation you want to perform.

- R also knows math functions like square root (`=sqrt()`)
- R commands will be shown alongside their resulting outputs, if entered in the R console.

R packages

What is an R package?

- Collection of functions, data and documentation that explains the package
- It adds useful and extra functionality not included in base R

Download and install a package

- `install.packages()`, e.g., `install.packages("tidyverse")`
- You only need to install a package once

Load the installed package

- `library()`, e.g., `library(tidyverse)`
- You need to load a package each time you open RStudio
 - o Makes all its functions available to use
 - o Coding rule: Load the packages at the start of your R session or script

Objects

R can store information as an object. This can be everything from a whole dataset to a certain text (=string variable)

- Use the assignment operator (`<-`) to assign some value to an object.
- `Object_name <- value`
 - o E.g. `x <- 15` or `y <- "Hi, welcome to the first programming lecture."`
- The created object will appear in your workspace window.
 - o Coding rule: Use informative names for files, variables, and functions.
- Assigning a new value to an existing object overwrites. The previous value.

R recognizes different classes of objects.

- Assigning object into classes allows R to perform class-specific operations
- Use `class()` to check the class of an object.
 - o Classes can vary from "numeric", "character" to "function" (=e.g. `sqrt`)

Vectors

A one-dimensional collection of elements, all of the same data type

- They enable efficient data manipulation

Create vectors

- `c()` function, which stands for "concatenate"
- Commas separate different element of a vector
 - o E.g. `c(TRUE, FALSE, True, NA)`

Types of vectors (`typeof()`)

- logical vector: FALSE, TRUE, and NA (e.g. `c(TRUE, FALSE, TRUE, NA)`)
- Numeric vector: Numbers (e.g. `c(1, 2, 3)`)
- Character vector: Strings (e.g. `c("NL", "BE", "DE")`)

Access specific elements of a vector ("indexing")

- By using `[]`
 - o E.g. `b[1]` in vector `b <- c(TRUE, FALSE, TRUE, NA)` gives TRUE.

Combine multiple vectors into one by simply putting the vectors you want to add into a new vector.

Avoiding errors

Encountering errors is part of the learning process, you'll need to learn how to fix these errors.

- Error: object not found
 - o R is case sensitive, so always check your spelling.
- Error: incomplete expression
 - o When you don't finish an expression you'll get the `+` prompt.
 - o Press ESC if you want to cancel the command being evaluated.

If you at some point don't know what to do for a specific function you can always use the help function.

- E.g. `help(class())`

Or use google or chatgpt for example.

Data import

The working directory is the location on your computer where R loads data from and saves data from.

- To display the current working directory: `getwd()`

You can set a working directory where to store all your data on the computer.

- Use: `setwd()`
- Relative paths are flexible if you move your project folder.
 - o `..`: parent directory (e.g., C:/Users)
 - o `.`: current directory (e.g., C:/Users/Programming)
- Absolute paths specify the full location (e.g. C:/Users/Programming/Week 1)
- Use `/` as path separator

Tidyverse

There are many ways to write code in R that produce the same results

- The tidyverse is like a dialect or syntax style of R programming
- It provides a more efficient and elegant way to work with data
- It is build for data science tasks: importing, cleaning, and transforming

To load files in R you'll need the readr package (=part of tidyverse). For different filetypes you'll need different functions

- Txt. Files: `read_delim()`
- Csv. Files: `read_csv()` (comma delimited) and `read_csv2()` (semi-colon delimited)
 - o The `col_names` argument is standard set at TRUE, this means that the first line of the file being read contains the variable names.
 - o If this is not the case, you can use `skip = n` to skip the first n lines
 - o Use `comment = "#"` to drop all lines that start with #

A **tibble** or **data.frame** object is a two-dimensional table with rows and columns, where each column can be of a different data type

- Use `view()` to view a table-like representation of *tibble* or *data.frame* objects.

Data export

Use the function `write_csv()` to save data as a CSV file

- E.g. `write_csv(name dataset, "Filename")`

- Use a new filename if you don't want to replace the original file

To view the first rows of data you can use: `head()`

- E.g. `head(name dataset)`
- If you only want to see the first n rows you can add it to the function
 - o E.g. `head(name dataset, n)`

To view the last rows of data you can use: `tail()`

- E.g. `tail(name dataset, n)`

Summarizing data characteristics

To get insights into each variable's distribution and characteristics use `summary()`

- This helps in understanding variable types, ranges, and missing values
 - o E.g. `summary(name dataset)`
- Check 1 specific column of the data by using `$`
 - o E.g. `summary(name dataset$name column)`
- Check multiple columns of the data by using `[]`
 - o E.g. `summary(name dataset [1:5])`
 - o This shows the summary of column 1 to 5.

Dataset structure

The `dim()` function checks the dimension of the dataset

- E.g. `dim(name dataset)`
- It returns 2 values
 - o The number of observations (rows)
 - o The number of variables (columns)

The `str()` function checks the structure of the dataset

- E.g. `str(name dataset)`
- It returns a summary of the dataset, such as:
 - o Number of observations (rows) and variables (columns)
 - o Data types of each column (e.g. integer, numeric, character, factor)
 - o The first few values in each column.

R markdown

R markdown is a plain text file format that combines code, text, and results in a single document. It support various output formats, such as PDF and Word, and it generated high-quality, fully reproducible reports for easy sharing and collaborating. To create a

new R Markdown document you do the following: File -> New File -> R Markdown. To generate reports in PDF, Word or HTML format you use the **knit** function in the toolbar.

An R Markdown document contains **3 types of content**

1. YAML header surrounded by: ---
 - o Specifies title, author, date and output format.
2. Chunks of R code surrounded by: ```
3. Text with simple text formatting (#=first level header, ##=second level header)

Run code in R Markdown:

- Insert a code chunk: Ctrl/Cmd+Alt+I
- Run current code line: Ctrl+Enter
- Run current code chunk: Ctrl+Shift+Enter

Chunk options

- Echo=FALSE -> show results, not code in finished line
- Eval=FALSE -> don't run the code
- Include=FALSE -> run the code, but code & result aren't shown
- Message=FALSE -> prevent messages in finished line
- Warning=FALSE -> prevent warning in finished file
- Results='hide' -> hide printed output
- Fig.show='hide' -> hide plots

Lecture 2

Tidying data

A large portion of data analysis time is spent cleaning and preparing data. Tidy datasets are easy to model, visualize and organized in a specific structure. But the main goal of tidying data is to store data values in a uniform way to facilitate analysis.

There are **3 interrelated rules** which makes a dataset tidy:

1. Each variable must have its own column.
2. Each observation must have its own row.
3. Each value must have its own cell.

All packages in the *tidyverse* are designed to work with tidy data. You can tidy up messy data with the *tidyr* package (is part of *tidyverse*).

Pivoting

Pivoting a table allows to change the form of your data without changing any of the values. You can pivot from wide to long: `pivot_longer()`, or from long to wide: `pivot_wider()`.

Long

i	j	Grade
1	1	8
1	2	4
2	1	6
2	2	7

Wide

i	Grade1	Grade2
1	8	4
2	6	7

Wide to long

- `Pivot_longer(data, cols, names_to="name", values_to="value")`
 - o **Data**: data frame
 - o **Cols**: columns to pivot into a longer format
 - o **Names_to**: name of the column created from the data stored in the column names
 - o **Values_to**: name of the column created from the data stored in cell values.

`Pivot_longer(table, c("2024","2025"),
names_to="Year", values_to="Grade")`

ID	2024	2025
1	8	4
2	6	7




ID	Year	Grade
1	2024	8
1	2025	4
2	2024	6
2	2025	7

Long to wide

- `Pivot_wider(data, names_from="name", values_from="value")`
 - o `Data`: data frame
 - o `Names_from`: name of the column to get the names of the output columns from.
 - o `Values_from`: name of the column to get the cell values from.

`Pivot_wider(table, names_from="Year", values_from="Grade")`

ID	Year	Grade
1	2024	8
1	2025	4
2	2024	6
2	2025	7



ID	2024	2025
1	8	4
2	6	7

Separating and uniting columns

Sometimes it can happen that there's multiple data in one cell (e.g. lastname/firstname) or if you want to combine two cells. Then you can use one of the two following codes:

- `Separate(data, col, into, sep)`
 - o `Data`: data frame
 - o `Col`: column to separate
 - o `Into`: names of new columns
 - o `Sep`: separator between columns
 - o E.g. `separate(table, Name, into=c("Lastname", "Firstname"), sep="/")`
- `Unite(data, col, columns, sep)`
 - o `Data`: data frame
 - o `Col`: name of the new column
 - o `Columns`: columns to unite
 - o `Sep`: separator between values
 - o E.g. `unite(table, col=Name, c("Lastname", "Firstname"), sep="/")`

Handling missing values

Datasets often have missing values (NA). If you ignore these missing values this can lead to biased analyses or broken visualizations.

To detect missing values, you can use 2 functions:

- `is.na()`

- This tests each element for missing values (Gives TRUE if it's missing)
- `sum(is.na())`
 - Counts the number of missing values.

You can drop the rows with missing values using `drop_na()`.

To replace missing values you can use `replace_na()`

Data visualization: Tables

Use the `stargazer` package to produce high-quality summary and regression tables. Tables can directly be exported to LaTeX, HTML or text formats. This is especially useful when working with R markdown.

Summary table

Descriptive statistics in text format, results can be seen in R console.

- e.g. `stargazer(as.data.frame(df), type="text", title="new title", digits = 1, out="table.txt")`

To ensure that the LaTeX table code is directly inserted into the document without being altered you have to put: `results = "asis"` into the chunk.

Data visualization: Graphs

In R you can create graphs with the `ggplot2` package (part of tidyverse).

To make a graph you use the following functions:

- Define the dataset to use in the graph
 - `ggplot()` -> e.g., `ggplot(data=df)`
- Specify which variables to map to the x and y axes
 - `aes()` -> e.g., `ggplot(data=df, mapping=aes(x=income, y=experience))`
- Define the plot to draw
 - `geom_function()` -> e.g., `ggplot(data=df, mapping=aes(x=income, y=experience)) + geom_point()`
- Save the graph to current working directory
 - `ggsave()` -> e.g., `ggsave(filename = "Figure 1.png")`

Common types of graphs:

- `geom_point()`: scatter plot
- `geom_histogram()`: Histogram
- `geom_line()`: lines/time series plot
- `geom_bar()`: bar plot
- `geom_boxplot()`: Box plot

- `geom_density()`: Density estimate
- `geom_smooth()`: fitted regression line

Example: Time series plot

The `geom_line()` function draws a connected line. You can also filter the data with the following function: `filter(data, condition)`

Advanced graphics

You can also add some aesthetics to the graph. The most common are:

- color -> color for point/lines
- fill -> fill color for areas
- shape -> symbol for points (square, round, triangle etc.)
- linetype -> type of line (solid, dashed, dotted etc.)
- size -> size of points/lines

The `labs()` function allows you to set labels for multiple plot elements all at once.

It is also possible to work with multiple layers in your code. In this way you can show a scatterplot and a regression line in the same graph. To do so you can just type the codes on top of each other.

Global vs local aesthetics

- Global
 - o Aesthetics are globally defined for the entire plot. Mappings for x, y and color apply to all layers in the plot.
 - o `ggplot(data=df, mapping=aes(x=income, y=experience, color=age)) +
 geom_point() +
 Geom_smooth()`
- Local
 - o Aesthetics are locally defined within each layer
 - o `ggplot(data=df) +
 geom_point(mapping=aes(x=income, y=experience, color=age)) +
 Geom_smooth(mapping=aes(x=income, y=experience))`

Lecture 3

Data transformation

Transform data with the *dplyr* package, which is part of *tidyverse*

Basic grammar of data transformation

- Operation on **rows**
 - o **filter()**: picks rows of a table based on their values
 - o **arrange()**: orders rows based on values of a column(s)
- Operation on **columns**
 - o **select()**: picks columns of a table based on their names
 - o **mutate()**: adds new columns that are functions of existing columns
- Operation on **groups**
 - o **summarize()**: aggregates multiple rows down to a single summary
 - o **group by()**: divides a table into groups

Filter

filter(data, ...): picks rows of a table based on their values that meet a certain condition.

- **data**: data frame
- **...**: Condition(s) that must be true to keep a row
- E.g. `storms_a <- filter(storms, wind >= 50)`

Logical and Boolean operators to use with filter()			
==	equal to	is.na()	is missing
!=	not equal to	!is.na()	is not missing
<	less than	%in%	in a set of specific values
>	greater than		or
<=	less than or equal	&	and
>=	greater than or equal		

Arrange

arrange(data, column1, column2, ...): orders rows based on values of a column

- **Data**: dataframe
- **Column1, column2**: Columns by which to arrange the data
- **desc()**: descending order (ascending is standard)
- E.g. `storms_b <- arrange(storms, wind)`

Select (same as filter function but now for rows, easy reminder select has a c of columns and Filter has a r of rows)

select(data, columns_to_select): picks columns of a table based on their names

- **Data**: data frame
- **Columns_to_select**: columns to keep
 - o Direct column names: **column names**, separated by commas
 - o Range of columns: **:** to select a range of columns
 - o Column exclusion: **-** to exclude columns

- E.g. `storms_c<-select(storms, -Date)` (Here you remove the column date)

Mutate

`mutate(data, ...)`: adds new columns that are functions of existing columns

- `Data`: data frame
- `...:` expressions that define the new columns (`new_column_name = expression`)
 - E.g. `storms_d<-mutate(storms, ratio=pressure/wind)`

If else

`if_else()`: create conditional variable

- `Condition`: logical vector
- `TRUE`: where condition is TRUE, matching values from TRUE
- `FALSE`: where condition is FALSE, matching values from FALSE
- E.g. `storms_e<-mutate(storms, red=if_else(wind>100, 1, 0))`

Summarize and group by

`Summarize(data, ...)`: aggregates multiple rows down to a single summary

- `Data`: data frame
- `...:` One or more expressions that calculate summary statistics (`new_column_name = summary_function(column_name)`)
 - Summary functions: `sum()`, `mean()`, `sd()`

`Group_by(data, ...)`: divides a table into groups

- `Data`: data frame
- `...:` columns to group by

The pipe operator

- `%>%`
- This function allows to combine multiple function without always having to repeat for example the data frame.
- E.g. (In this case you don't have to repeat "table1" all the time in the functions)


```
table1 %>%
  Group_by(Year) %>%
  Summarize(Mean_Case=mean(Case))
```

Ungroup

The function `ungroup()` removes existing grouping structure created by the `group_by()` functions

- If you don't use the `ungroup()` function averages are calculated **within** each group, when you do use it, averages are calculated **across** all rows in the dataset.

Row-wise

Row-wise operations require grouping of a row

- Built-in functions: `rowSums()`, `rowMeans()`
- To select which columns to include use: `across()`
- To handle missing values add: `na.rm=TRUE`
 - o E.g. `total_score=rowSums(across(c(Column1, column2)), na.rm=TRUE)`

Joining data

Combine data with the *dplyr* package, which is part of *tidyverse*

Basic grammar for combining data

Operation on tables

- Append columns (datasets have the same number of rows, but different columns)
 - o `Bind_cols(x,y)`: combines tables by adding columns side-by-side
- Append rows (datasets have same columns but different rows)
 - o `Bind_rows(x,y)`: combines tables by stacking rows on top of each other.
- Mutating joins
 - o `inner_join(x,y)`: keeps observations that occur in both x and y
 - o `left_join(x,y)`: keeps all observations in x
 - o `right_join(x,y)`: keeps all observations in y
 - o `Full_join(x,y)`: keeps all observations in either x or y
- Filtering joins
 - o `semi_join(x,y)`: keeps all observations in x that have a match in y (but you'll only keep the observations from x; no new column is added)
 - o `anti_join(x,y)`: drops all observations in x that have a match in y (You'll only get the observations that aren't in y)

Keys: variables used to connect a pair of data frames in a join

- **Primary key:** variable(s) that uniquely identifies each observation.
- **Foreign key:** variable(s) that corresponds to a primary key in another table.

Duplicate values

Duplicate values can appear when merging, importing, or appending data

- Identical values across entire rows or within specific columns

Detect duplicates: `duplicated()`

- You'll get a logical vector: TRUE for duplicates

Count duplicates: `sum(duplicated())`

- It will sum up the TRUE values.

Locate duplicates: `which(duplicated())`

- It will give the row number of the duplicate.

Remove duplicates: `distinct()`

Dates in R

To work with dates and times in R, you'll need the *lubridate* package, which is part of the *tidyverse* package.

Depending on how the date is written, you can convert character or numeric date values into date objects.

- **Y** -> year
- **M** -> month
- **D** -> day

Different functions

- `ymd(20221101)` -> "2022-11-01"
- `dmy(01112022)` -> "2022-11-01"
- `mdy("November 1st, 2022")` -> "2022-11-01"

It is also possible to calculate how much time elapsed between two dates

- `entry_date <- ymd("2022-06-01")`
`implementation_date <- ymd("2024-01-15")`
`implementation_date - entry_date`

Time difference of 593 days

- You can also use `as.numeric(x)`, this converts X into a number
 - o `as.numeric(implementation_date - entry_date)`
[1] 593

Get individual components from a date variable by using `$` (output is numeric!)

- Get the year: `df$year <- year(df$date)`
- Get the month: `df$month <- month(df$date)`
- Get the day: `df$day <- mday(df$date)`

Create a date variable from individual components with `make_date()`

- Input = numeric
- Output = date data type
 - o E.g. `df$date_new <- make_date(df$year, df$month, df$day)`

Lecture 4

Regression analysis

A regression analysis is a statistical method for analyzing a relationship between two or more variables. In other words, one variable can be predicted or explained by using information on the other variables.

- Idea: fit a function that closely represents the trend in the data.

The most common form of regression analysis is the **linear regression**. This fits a straight line to the data.

- $Y = \alpha + \beta X + \varepsilon$
 - o Y: Dependent, response or outcome variable
 - o X: Independent, explanatory or predictor variable
 - o α : Intercept
 - o β : slope coefficient
 - o ε : error term
- It uses the ordinary least squares

In R you can also do a linear regression by using the `lm(formula, data)` function.

- **Formula**: `Y ~ X (+...)`
- **Data**: data frame

You can do both a **univariate** and a **multivariate** linear regression

- Univariate linear regression (=simple linear regression)
 - o E.g. `unireg <- lm(sales ~ advertising, data=sales)`
- Multivariate linear regression
 - o E.g. `multireg <- lm(sales ~ advertising + rd, data = sales)`

The `summary()` function shows then the regression output.

```
Call:
lm(formula = distance ~ avg_quality + fr_orig + de_orig + sc_orig,
    data = country_data)

Residuals:
    Min       1Q   Median       3Q      Max
-24.1910  -2.8133   0.0713   3.4919   9.3159

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)    8.315      2.246   3.702 0.000251 ***
avg_quality   -8.620      1.096  -7.866 5.47e-14 ***
fr_orig1      -2.023      1.495  -1.353 0.176948
de_orig1       2.124      1.517   1.401 0.162273
sc_orig1       1.801      1.580   1.139 0.255339
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.703 on 325 degrees of freedom
Multiple R-squared:  0.2846,    Adjusted R-squared:  0.2758
F-statistic: 32.33 on 4 and 325 DF,  p-value: < 2.2e-16
```

Call: model specification

Residuals: The difference between the observed and predicted values (ideally this should be symmetrically distributed around the line)

Coefficients: Significance testing (standard error, t-value and p-value)

Residual standard error: std. dev. of residuals

The function `stargazer()` can make a table out of the regression output.

- E.g. `stargazer(unireg, multireg, type="text", title="sales")`

It is also possible that you use categorical variables. In this case you can use a **logistic regression**. This fits a regression curve when Y is a categorical variable. These variables have a fixed and known set of possible values.

- Binary logistic regression (A or B)
- Multinomial logistic regression (A,B, C)
- Ordinal logistic regression (A<B<C)

In R the `glm(formula, family, data)` creates a logistic regression model.

- **Formula:** Y~X (+...)
- **Family:** probability distribution (for logistic regression: `family=binomial(link="logit")`)
- **Data:** data frame

Generalized linear models allow for non-normally distributed dependent variables.

Before you can execute the logistic regression, you will have to do some preparatory steps. It is important to transform categorical variables to factors. You can do this by using `as.factor()`.

The estimated coefficients are in log odds scale $\rightarrow \log\left(\frac{p}{1-p}\right)$

Fixed effects

Fixed effects control for unobservable confounding factors that are constant within categories. It is a variable that helps rule out alternative explanations for the

relationship between the outcome and the explanatory variable of interest. What is it in 3 points:

- Fixed effects are dummy variables for each category
- Used as covariates in regressions to account for unobservable factors
- Add category-specific constants

3 types of fixed effects

1. Unit fixed effects

- o Account for fixed differences across units (e.g. industries, countries)
- o Basically, characteristics that vary across units but do not vary over time.
- o By including state fixed effects, you control for these time-invariant characteristics.

2. Time fixed effects

- o Account for common time trends
- o Characteristics that vary over time, but do not vary across units.
- o By including year fixed effects, you control for these common time-varying factors.

3. Two-way fixed effects

- o Regression model including both unit and time fixed effects
- o Two-way fixed effects design accounts for confounding factors that are either:
 - Constant over time but vary across units.
 - Constant across units but vary over time.
- o Two-way fixed effects allow you to isolate the within-state relationship between the dependent and independent variable.

In R you can include fixed effects as well. For this you'll need the lfe package.

Basic syntax:

felm(regression formula, FE dimension, instrumental variable, SE clustering dimension, data frame)

Panel data

3 types

1. Cross-sectional data
 - o One point in time ($t=1$), many observation units ($i=1, \dots, n$)
2. Time-series data
 - o Many time periods ($t=1, \dots, T$), one observation unit ($i=1$)
3. Panel data

- Many observation units ($i=1, \dots, n$) over time ($t=1, \dots, T$)

Loops

Loops is a tool for reducing duplication of effort by automatically repeating a set of operations several times. There are 3 types of loops:

1. **for** loop
2. **while** loop
3. Repeat loop

e.g. merging 50 data files into one dataset.

for loop

Repeat a task a defined number of times.

Basic syntax

- `for(element in vector){
 expression
}`
- e.g. `for(i in 1:5) {
 print(i)
}`

This line of code will print `i` that is each value in the vector, so it will give you the numbers 1, 2, 3, 4 and 5

while loop

It repeats a task until a specific condition is met.

Basic syntax:

- `while(logical condition) {
 expression
}`
- e.g. `i <- 0
while(i<=4){
 i <- i + 1
 print(i)
}`

This will give the output 1, 2, 3, 4 and 5.

- If you don't add the condition `i <- i + 1`, the loop will run indefinitely because `i` will always be 0 and thus meet the condition.

Functions

R is full of **built-in functions**. These are predefined functions that are available in R to perform common tasks or operations. R does also allow you to write your own function (= **user-defined function**). This function is a set of statements organized together to perform a specific task.

Basic syntax:

- `f <- function(argument(s)) {`
 statements
 }
- This creates a function with name `f` which takes certain arguments and executes the following statements.
- E.g. create a function that adds two numbers:
 `Add_num <- function(a,b) {`
 `Result <- a+b`
 `Print(result)`
 }
- Call the function:
 `Sum <- add_num(3, 4)`
 This will give the result 7

Lecture 5

String manipulation

Textual analysis is the process by which information is extracted and classified from unstructured text data. There are 2 most widely used techniques in textual analysis:

1. Sentiment analysis: Identify the underlying sentiment or text
2. Topic detection: Group themes.

It is possible to manipulate strings (text stored in a variable) with the *stringr* package (part of *tidyverse*). All the functions of the *stringr* package start with `str_`

There are 11 important `str_` functions:

- **`str_c()`**: merge several strings into a single string
 - o E.g. `str_c("a", "b", "c")` will get the output "abc"
- **`str_split()`**: split a string up into pieces
 - o E.g. `str_split("split me", pattern=" ")` will get the output "split" "me"

Working with individual characters within a string

- **`str_length()`**: determine a string's number of characters

- E.g. `str_length(c("p", "programming in R", NA))` will get the output 1 18 NA
 - Spaces and punctuation are counted as characters, missing values not.
- **Str_sub()**: extract characters of a string (by position)
 - E.g. `y <- c("Orange", "Grape", "Apple")`
`Str_sub(y, 1, 3)` will get the output "Ora" "Gra" "App"
`Str_sub(y, -1, -3)` will get the output "nge" "ape" "ple"

Working with string matches (`x <- c("why", "video", "cross", "extra", "deal", "authority")`)

- **Str_subset()**: return all elements of string where there's at least one match to pattern
 - E.g. `str_subset(x, "[aeiou]")`
 Will give the output: "video" "cross" "extra" "deal" "authority"
- **Str_detect()**: return a logical vector (TRUE) for each element of string that matches the pattern.
 - E.g. `str_detect(x, "[aeiou]")`
 Will give the output: FALSE TRUE TRUE TRUE TRUE TRUE
- **Str_count()**: count the number of times pattern is found within each element of string.
 - E.g. `str_count(x, "[aeiou]")`
 Will give the output: 0 3 1 2 2 4
- **Str_extract()**: extract the first complete match from each string
 - E.g. `str_extract(x, "[aeiou]")`
 Will give the output: NA "i" "o" "e" "e" "a"
- **Str_extract_all()**: extract all matches from each string
- **Str_replace()**: replace the first match
 - E.g. `str_replace(x, "[aeiou]", "?")`
 Will give the output: "why" "v?ideo" "cr?ss" "?xtra" "d?al" "?uthority"
- **Str_replace_all()**: replace all matches
 - E.g. `str_replace_all(x, "[aeiou]", "?")`
 Will give the output: "why" "v?d???" "cr?ss" "?xtr?" "d??" "??th?r?ty"
 - Or perform multiple replacements (`x <- c("1 house", "2 cars")`)
 - e.g. `str_replace_all(x, c("1"="one", "2"="two"))`
 will give the output: "one house" "two cars"

Regular expressions

What are regular expressions?

- A sequence of characters defining a text pattern.
- It is useful for finding, extracting, or replacing text.

It is most used for:

- Detecting valid phone numbers
- Matching email addresses
- Checking password rules
- Cleaning company names

Literal characters -> Characters that match themselves (every letter of the English alphabet or number)

Str_view(): Print the string and see how a pattern matches.

- E.g. `this_flower <- "This flower is mine. I bought this flower from the flowershop"`
`Str_view(this_flower, "flower")`
Will give the output: This <flower> is mine. I bought this <flower> from the flowershop
- It is also possible to match parts of words, for example only "flow"

Metacharacters -> Characters that have a special meaning (`.\^$[]- *+?`)

Wild metacharacters -> Matches every character except for a new line

- `.` | any character
- E.g. `x <- c("apple", "banana", "pear")`
`Str_view(x, ".a.")` (<- the `.` means any character before or after "a")
Will give the output: <ban>ana, p<ear>

To escape metacharacters you can use: `\\`

Anchors -> Match the position of a pattern. `^` means the start of a string, `$` means the end of a string.

- E.g. `str_view(x, "^a")` will give: <a>pple
- `Str_view(x, "a$")` will give: banan<a>

Character sets -> Match a set or range of characters.

- `[]` defines a character set (`[abc]` means a, b or c, `[a-z]` means a range)
- `[^]` means anything but (`[^abc]` means anything but a, b or c)
- Metacharacters inside a character set are already escaped, so you don't need `\\` (`"p[ae.io]n"` means pan, pen, p.n, pin, pon)
 - o `.` inside a character set is a literal dot, not everything around it

Character classes -> match a certain class of characters

Character	Matches	Same as
-----------	---------	---------

\\d	Any digit	[0-9]
\\D	Any non-digit	[^0-9] (=anything but 0-9)
\\w	Any word character	[a-zA-Z0-9_] -> matches letters, digits, and underscores
\\W	Not a word character	[^a-zA-Z0-9_]
\\s	Whitespace	
\\S	Any non-whitespace character	
\\b	Word boundary	

Quantifiers -> These specify how many times an element should be matched

*	Zero or more matches
+	One or more matches
?	Zero or one match
{m}	Exactly m matches
{m, }	M or more matches
{m, n}	Between m and n matches

- E.g. `names <- c("Dr. Clark", "Carol", "Sameer", "George Jr.")`
`Str_view(names, "^[A-Z][a-z]*$")` -> *means string consists of a single uppercase letter followed by zero or more lowercase letters.*
Will give the output: `<Carol>`, `<Sameer>`

Look arounds -> Match only if specified pattern can(not) be matched before (after)

?=	Followed by	Positive look ahead
?!	Not followed by	Negative look ahead
?<=	Preceded by	Positive look behind
?<!	Not preceded by	Negative look behind

- E.g. `weights <- "55lbs, 8kg, 32, 14lbs"`
`Str_view(weights, "[0-9]+(?=lbs)")` -> *means any number followed by lbs*
Will give: `<55>lbs`, 8kg, 32, `<14>lbs`

Lecture 6

Recap

Tidyverse is a collection of R packages for data science:

- Week 1: Data import -> *readr*
- Week 2: Data tidying and data visualization -> *tidyr* and *ggplot2*
- Week 3: Data transformation -> *dplyr* and *lubridate*
- Week 5: strings -> *stringr*

Other R packages for data science:

- Week 2: data tidying and data visualization -> *stargazer*
- Week 4: Regression analysis, loops, & functions -> *lfe*

Tidy text format

A tidy text format is a data frame with one-token-per-row. A token is a meaningful unit of text (such as a word, n-gram, sentence, or paragraph).

But why should tidy text data? Tidy text format conforms to tidy data principles and can be manipulated with a set of consistent tools. It allows manipulation with popular packages such as *dplyr*, *tidyr* and *ggplot2*.

- In order to convert text to a tidy format you have to install & load the *tidytext* package.

Example

Create a character vector

- `Text <- c("Mr. and Mrs. Dursley, of number four, Privet Drive, were",
"proud to say that they were perfectly normal, thank",
"you very much.")`

Put the character vector into a data frame

- `Text_df <- tibble(line=1:5, text=text)`

Convert the data frame into one-token-per-row (=tokenization). This transforms the text to a tidy data structure.

- `Unnest_tokens(output, input)`: split a column into tokens
 - o **Output**: output column to be created
 - o **Input**: input column to be split
 - o It retains other columns, removes punctuation and converts tokens to lowercase.
- E.g. `text_df %>%
Unnest_tokens(word, text)`

Word Frequency and Word Cloud

A word frequency analysis is a common text mining technique that involves counting how often each word appears in a collection of text. It helps identify the most frequently occurring words. This can be useful for understanding the importance and prevalence of words within the text. This can gain valuable insights into the structure and content of textual information.

Text mining requires a lot of data cleaning. If you have formatted a data frame in one-word-per-row then you notice that not all words are useful for analysis. There are the standard stop words (most commonly used words like “a”, “the”, “is”, “of”).

- `Stop_words`: *tidytext* list of stopwords
 - o English stop words from three lexicons as a data frame.
- `Get_stopwords()`: get a specific stop word lexicon
- To remove standard stop words you can use: `anti_join(stop_words)`
- `Count()`: count the unique values of one or more variables
 - o Always think about the results, because it could be possible that a word has a double meaning (e.g. miss).

Effective word frequency analysis goes beyond counting, it requires identifying the meaningful terms and interpreting the context in which they appear. Counting words is simple, but identifying the relevant ones can be tricky:

- Simple keywords: “covid”, “risk”
- Pattern-based sequence classification: “supply chain disruption”
- Advanced techniques: sophisticated keywords and embeddings

You can also leverage context from large textual data, for example:

- For climate change: What specific concerns do firms have, such as regulatory impacts or rising operational costs?
- For Brexit: What specific worries do firms face, such as fears of labor shortages or currency fluctuations?

A common way to visualize a word frequency distribution is a **word cloud**. Words that are more frequently used, appear in larger font.

- You can create a word cloud visualization with the *wordcloud* or *wordcloud2* package.

Sentiment Analysis

A sentiment analysis is a text mining technique to determine the emotional tone of text (positive, negative or neutral). It uses the emotional intent of words to infer the sentiment of text. Sentiment analysis is useful for understanding the underlying mood or opinions.

Sentiment lexicons

- **AFINN**: assigns words with a score that runs between -5 and 5

- **Bing**: categorizes words in a binary fashion (yes/no) into positive and negative categories.
- **Nrc**: categorizes words in a binary fashion (yes/no) into categories of positive, negative, anger, anticipation, disgust, fear, joy, sadness, surprise, and trust.

Dictionary/sentiment word list in financial contexts (domain-specific)

- Loughran-McDonald: seven sentiment identifiers
 - o Negative, positive, uncertainty, litigious, strong modal, weak modal, constraining.

Dictionary-based methods find the total sentiment of a piece of text by adding up the individual sentiment scores for each word in the text.

What positive words exist in *bing* dictionary?

- **Get_sentiments()**: get specific sentiment lexicons.

Challenges & considerations

It can be difficult to handle negations. Simple word count may misinterpret sentiment.

- For example: "I am happy" or "I am not happy"

Figurative language and context (e.g. metaphors, irony, jokes). Sentiment analysis struggles with figurative language.

- For example: "He's such a genius!" (literal: positive; sarcastic: negative)