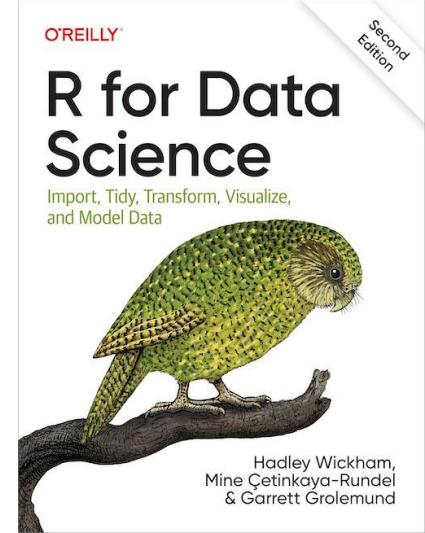**R I**

# Programmierkurs 2 Data Science WS23/24

Leonard Traeger
M. Sc. Information Systems
leonard.traeger@fh-dortmund.de

# Disclaimer

Slides are mainly based on

- https://r4ds.hadley.nz/

- https://www.phonetik.uni-muenchen.de/~jmh/lehre/basic_r/_book/index.html

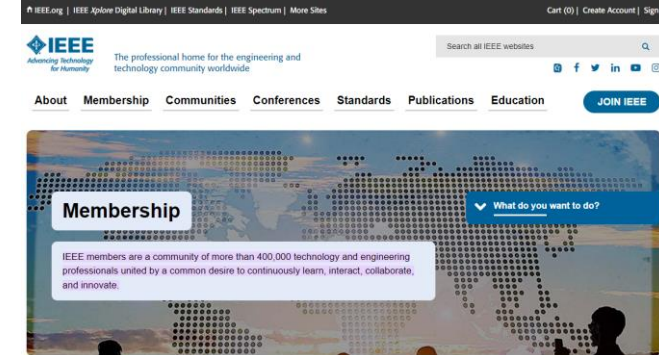→ Find everything you need to know there!


Official R cheat sheet:

- https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf

Data Transformation with dplyr:

- https://raw.githubusercontent.com/rstudio/cheatsheets/master/data-transformation.pdf

Technology
Arts Sciences
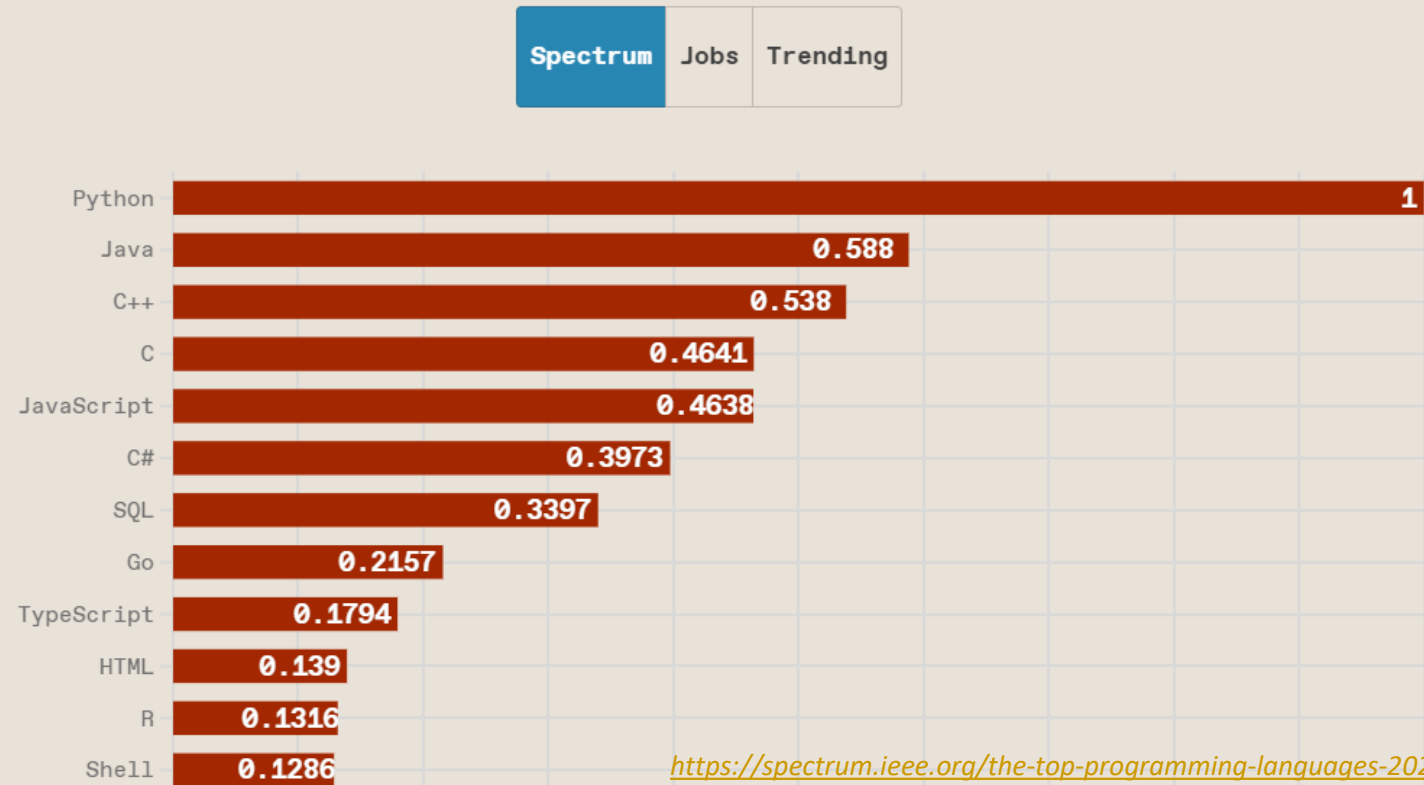**TH Köln**

# Learning Goals R I

- **Explain** your personal preference of R and Python as a programming language given by giving a comparitive coding example.

- **List** the development stack and its components for programming in R.

- Install and Import libraries and **use** R as a calculator.

- **Create** variables, vectors, matricies, and simple scatter plots.

- Import tabular files as DataFrames and **apply** exploration, filtering, slicing, selection, mutation, and renaming opreations using the pipe syntax.

Technology
Arts Sciences
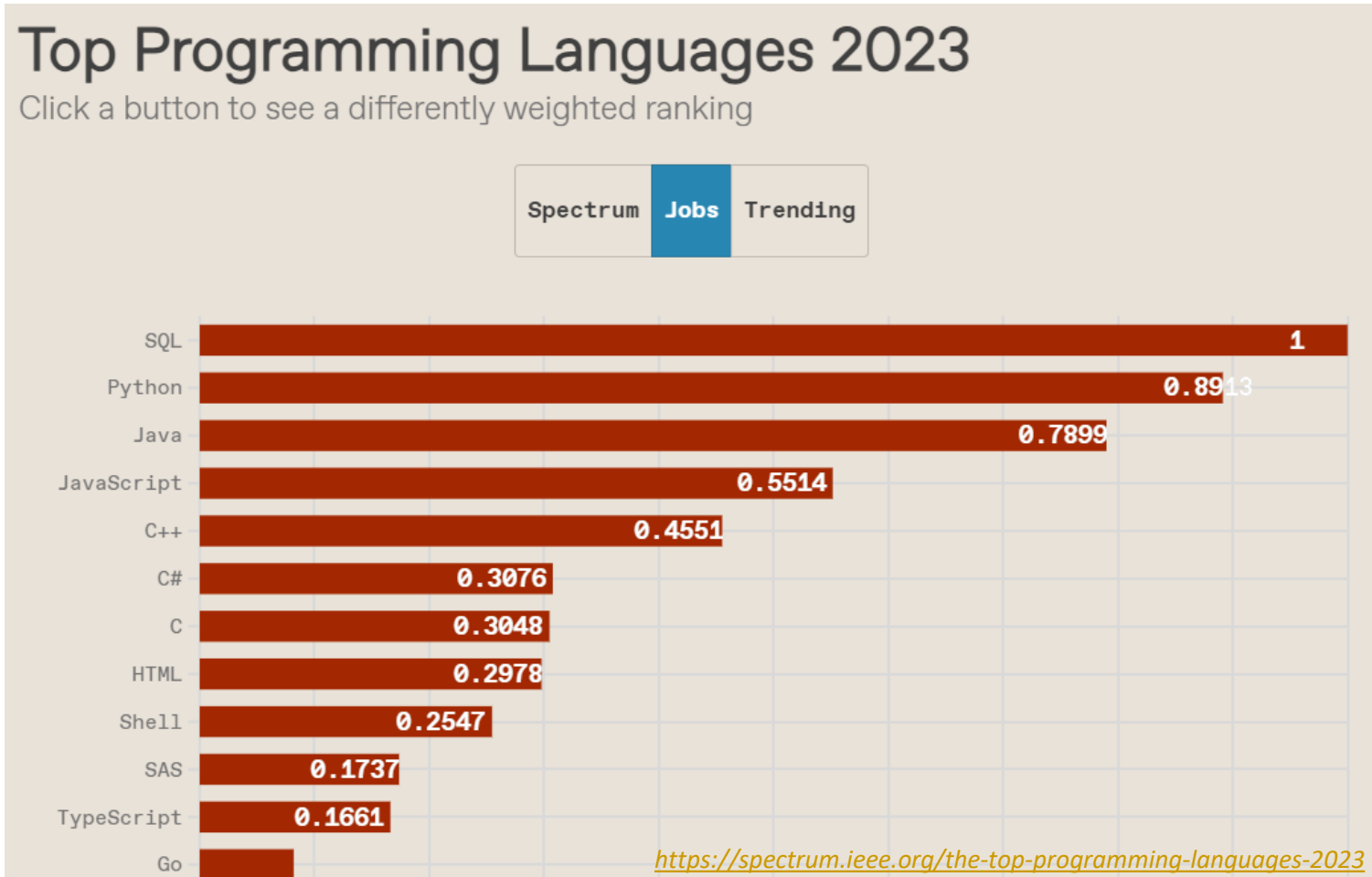TH Köln

# Top Programming Languages



Top Programming Languages 2023
Click a button to see a differently weighted ranking

Spectrum | Jobs | Trending

| | |
|---|---|
| Python | 1 |
| Java | 0.588 |
| C++ | 0.538 |
| C | 0.4641 |
| JavaScript | 0.4638 |
| C# | 0.3973 |
| SQL | 0.3397 |
| Go | 0.2157 |
| TypeScript | 0.1794 |
| HTML | 0.139 |
| R | 0.1316 |
| Shell | 0.1286 |

https://spectrum.ieee.org/the-top-programming-languages-2023

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Top Programming Languages



Top Programming Languages 2023
Click a button to see a differently weighted ranking

Spectrum | **Jobs** | Trending

| Language | Score |
|---|---|
| SQL | 1 |
| Python | 0.8913 |
| Java | 0.7899 |
| JavaScript | 0.5514 |
| C++ | 0.4551 |
| C# | 0.3076 |
| C | 0.3048 |
| HTML | 0.2978 |
| Shell | 0.2547 |
| SAS | 0.1737 |
| TypeScript | 0.1661 |
| Go | |

https://spectrum.ieee.org/the-top-programming-languages-2023

Programmierkurs 2 Data Science: R I

**Technology
Arts Sciences
TH Köln**

# Top Programming Languages



Top Programming Languages 2023
Click a button to see a differently weighted ranking

Spectrum | Jobs | **Trending**
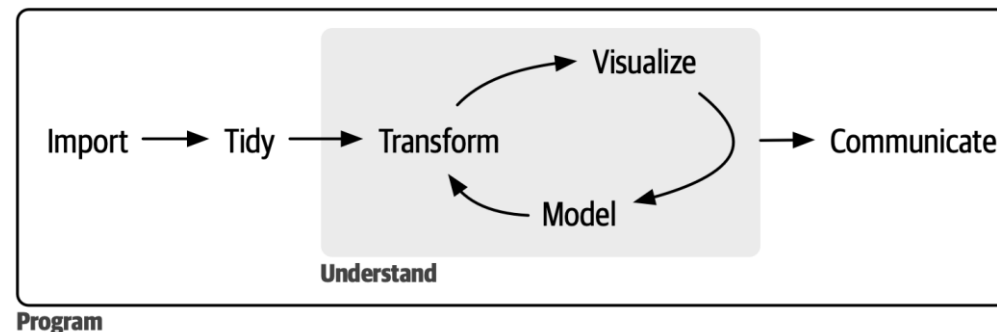
| Language | Value |
|---|---|
| Python | 1 |
| Java | 0.7358 |
| JavaScript | 0.554 |
| C++ | 0.4599 |
| SQL | 0.44?? |
| C# | 0.4187 |
| C | 0.3497 |
| HTML | 0.202 |
| Go | 0.1921 |
| TypeScript | 0.1882 |
| Shell | 0.142 |
| PHP | 0.1144 |
| R | 0.108 |

*https://spectrum.ieee.org/the-top-programming-languages-2023*

Programmierkurs 2 Data Science: R I

**Technology**
**Arts Sciences**
**TH Köln**

# Why R?

- Built to demonstrate the results of statistical analysis quickly.

- Suited for **statistical learning**.

- High level language; used by "non-techy" engineers and scientists.

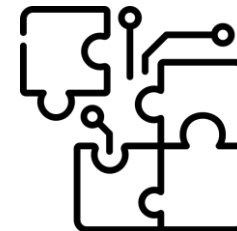- Open source, fast growing ecosystem with packages for almost everything in DS:
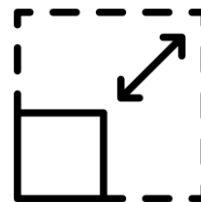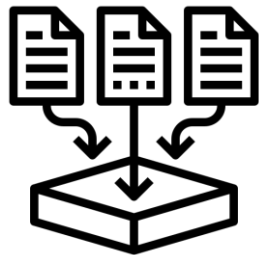


R for Data Science (e2) by Wickham, Çetinkaya-Rundel, and Grolemund

Programmierkurs 2 Data Science: R I

**Technology Arts Sciences TH Köln**

# Python versus R

|  | Python | R |
|---|:---:|:---:|
| Do you have programming experience? | ☺ | |
| Do you care about visualization and graphics? | | ☺ |
| Do you want to apply Statistical Models? | | ☺ |
| Do you want to apply Machine Learning? | ☺ | |
| What do your colleagues, peers, advisors, industry-area use? | ☺ | ☺ |

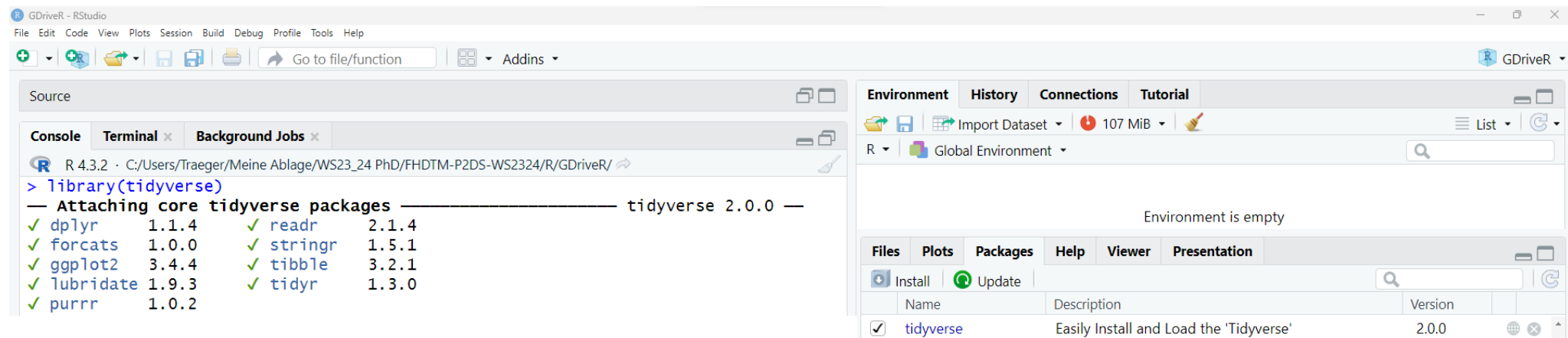Various decision factors: *data collection, libraries, scale, integration, and many more…*

Technology
Arts Sciences
TH Köln

# What you will need

**R**

**RStudio**: IDE for R programming

**Tidyverse**: a collection of R packages

- `install.packages("tidyverse")`

- `library("tidyverse")`

Technology
Arts Sciences
TH Köln

# Project in RStudio

File > New Project > Choose Directory



Advantages

- Restores the state of work where you left off.

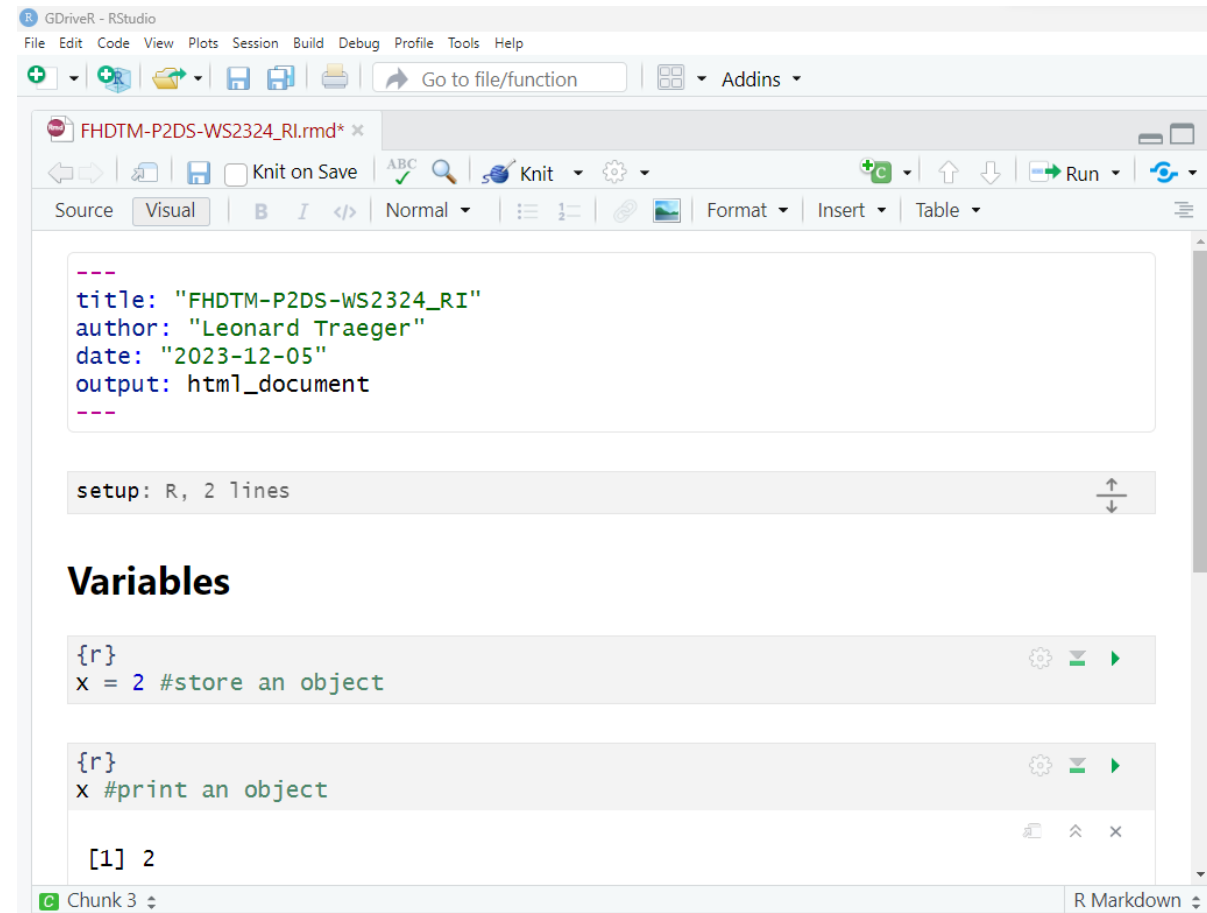- Files that you save during the course can be easily opened via the panel.

Programmierkurs 2 Data Science: R I

**Technology
Arts Sciences
TH Köln**

# R Markdown

File > New File > R Markdown > …
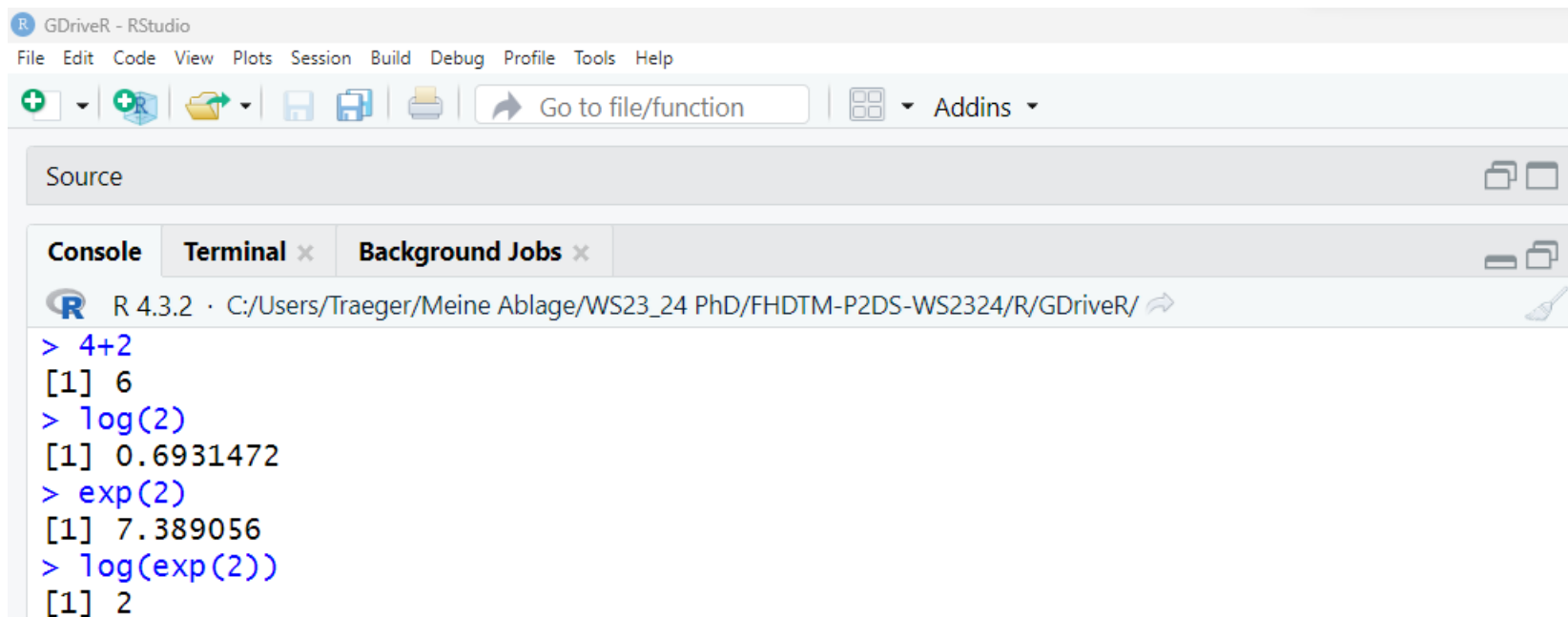
- Choose HTML type
- Store file as .Rmd

Text annotations

- # headline
- **bold**
- *italics*
- `code`



R Markdown documents has become established for the creation of report material (similar to our work in Python scripts).

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# RStudio Console

- The console is the direct connection between R and the computer that performs the calculations.

- You can use R as a simple calculator:



Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Numeric and string objects

R is a dynamically typed language.

```
> x = 2 #store an object
> x #print an object
[1] 2
> (y = 42) #store and print an object
[1] 42
> z = "Hello" #store a string object
> z
[1] "Hello"

> i <- 4 #object assignment via arrow operator
> 2 -> j #works also in the other way
> i
[1] 4
> j
[1] 2
```

R will ignore any text after # for that line.

Programmierkurs 2 Data Science: R I

**Technology Arts Sciences**
**TH Köln**

# Objects and Class

To find out which object class a variable has, use the `class(variable)` function.
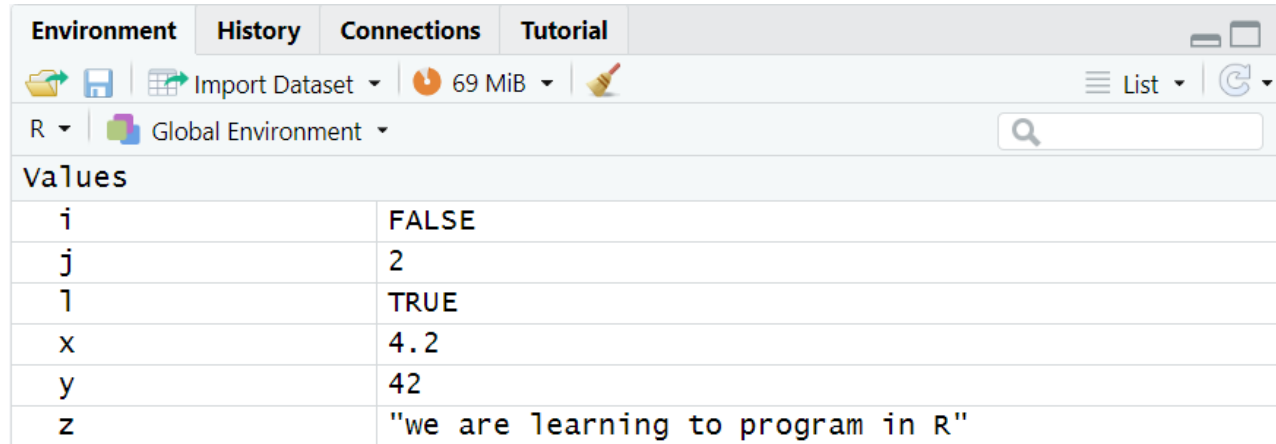
```
> x <- 4.2 #double
> z <- "we are learning to program in R"
> l = TRUE
> i = F #short for False
> class(x)
[1] "numeric"
> class(z)
[1] "character"
> class(l)
[1] "logical"
```

Programmierkurs 2 Data Science: R I

**Technology Arts Sciences**
**TH Köln**

# Environment Variables

- List all environment variables with the `ls()` function.

- Remove via `rm(variable)` function.

```
> ls()
[1] "i" "j" "l" "x" "z"
> rm(z)
> ls()
[1] "i" "j" "l" "x"
```

You can also see your environment variables in the top right window.

| Environment | History | Connections | Tutorial | |
|---|---|---|---|---|
| Import Dataset ▾ | 69 MiB ▾ | | | List ▾ |
| R ▾ | Global Environment ▾ | | | |

**Values**

| | |
|---|---|
| i | FALSE |
| j | 2 |
| l | TRUE |
| x | 4.2 |
| y | 42 |
| z | "we are learning to program in R" |

Programmierkurs 2 Data Science: R I

**Technology**
**Arts Sciences**
**TH Köln**

# Logical Operators

| | |
|---|---|
| `a < b` | Less than |
| `a > b` | Greater than |
| `a <= b` | Less equal than |
| `a >= b` | Greater equal than |
| `a == b` | Equal |
| `a != b` | Not equal |
| `!a` | Not |
| `a | b` | a OR b |
| `a & b` | a AND b |
| `isTRUE(a)` | Check whether a is TRUE |
| `a %in% c` | Check whether a's value is in a vector c |

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Vectors

Function `c()` (**concatenate**) creates a vector, a data structure with several elements.

If the elements belong to **different classes** (strings, booleans, numerics), the elements are **converted** to the **same type silently**, i.e. without a warning message.

```
> answer_to_everything = c("Bezos", "Zuckerberg", "Musk", 42)
> answer_to_everything[0] #0'th element stores the vector type
character(0)
> answer_to_everything[4]
[1] "42"
> answer_to_everything[2:3]
[1] "Zuckerberg" "Musk"
> numbers = c(0, 1, 2, 3, TRUE, FALSE)
> numbers
[1] 0 1 2 3 1 0
```

Programmierkurs 2 Data Science: R I
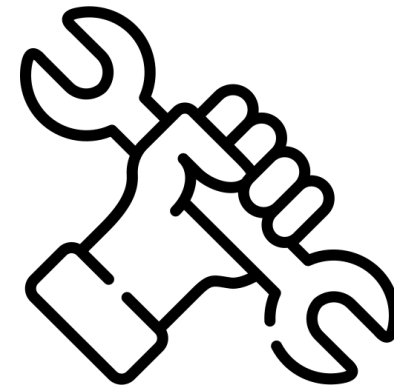
Technology
Arts Sciences
TH Köln

# Vectors (cont.)

Use the `seq(from, to, by)` function to create regular sequences of numbers:

```
> 1:5
[1] 1 2 3 4 5
> seq(from=1, to=5)
[1] 1 2 3 4 5
> seq(1, 5) #argument names are optional in R functions
[1] 1 2 3 4 5
> seq(1,5,by=0.5) #in intervals by 0.5
[1] 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0
```

Use the `rep()` function to create sequences with repeating values:

```
> rep(42, times=3)
[1] 42 42 42
> rep(c("Bezos","Musk"), times=3)
[1] "Bezos" "Musk"  "Bezos" "Musk"  "Bezos" "Musk"
```

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Training #1

1. Create a new project in Rstudio for this class.

2. Create two variables x = ''Hello'' and y = ''World'' and use them to print your first ''Hello World'' in R.
   Hint: use the `cat(a, b, ...)` function to print multiple variables on one line.

3. Return a vector with interchangble ''R'' and ''Python'' elements with the length 100.

You can access

- a virtual RStudio environment here: https://jup.labs.inf.fh-dortmund.de/

- RStudio over Anaconda

- Install R https://ftp.fau.de/cran/ and
  Rstudio https://posit.co/download/rstudio-desktop/#download

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Vectors (cont.)

You can apply basic **arithmetic operations** and arithmetic **functions** to numeric vectors:

```
> us_presidents_heights = c(189, 170, 189, 163, 183, 171, 185, 168, 173, 183,
+                           173, 173, 175, 178, 183, 193, 178, 173, 174, 183,
+                           183, 168, 170, 178, 182, 180, 183, 178, 182, 188,
+                           175, 179, 183, 193, 182, 183, 177, 185, 188, 188,
+                           182, 185, 191, 182)
> us_presidents_heights + 5
 [1] 194 175 194 168 188 176 190 173 178 188 178 178 180 183 188 198 183 178
[19] 179 188 188 173 175 183 187 185 188 183 187 193 180 184 188 198 187 188
[37] 182 190 193 193 187 190 196 187
> #performs single calculation using vectors
```

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Vectors (cont.)

You can apply basic **arithmetic operations** and arithmetic **functions** to numeric vectors:

```
> sqrt(us_presidents_heights)
 [1] 13.74773 13.03840 13.74773 12.76715 13.52775 13.07670 13.60147 12.96148
 [9] 13.15295 13.52775 13.15295 13.15295 13.22876 13.34166 13.52775 13.89244
[17] 13.34166 13.15295 13.19091 13.52775 13.52775 12.96148 13.03840 13.34166
[25] 13.49074 13.41641 13.52775 13.34166 13.49074 13.71131 13.22876 13.37909
[33] 13.52775 13.89244 13.49074 13.52775 13.30413 13.60147 13.71131 13.71131
[41] 13.49074 13.60147 13.82027 13.49074
> log(us_presidents_heights)
 [1] 5.241747 5.135798 5.241747 5.093750 5.209486 5.141664 5.220356 5.123964
 [9] 5.153292 5.209486 5.153292 5.153292 5.164786 5.181784 5.209486 5.262690
[17] 5.181784 5.153292 5.159055 5.209486 5.209486 5.123964 5.135798 5.181784
[25] 5.204007 5.192957 5.209486 5.181784 5.204007 5.236442 5.164786 5.187386
[33] 5.209486 5.262690 5.204007 5.209486 5.176150 5.220356 5.236442 5.236442
[41] 5.204007 5.220356 5.252273 5.204007
```

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Vector Aggregation

You can also describe numerical vectors with **aggregate** functions:

```
> length(us_presidents_heights)
[1] 44
> sum(us_presidents_heights)
[1] 7922
> mean(us_presidents_heights)
[1] 180.0455
> var(us_presidents_heights)
[1] 49.90486
```

Technology
Arts Sciences
TH Köln

# Vector Aggregation (cont.)

Use `unique(vector)` to display the unique items of a vector.

Use `table(vector)` to display a list of unique items and frequency of occurrence.

```
> unique(us_presidents_heights)
 [1] 189 170 163 183 171 185 168 173 175 178 193 174 182 180 188 179 177 191
> table(us_presidents_heights)
us_presidents_heights
163 168 170 171 173 174 175 177 178 179 180 182 183 185 188 189 191 193
  1   2   2   1   4   1   2   1   4   1   1   5   8   3   3   2   1   2
```

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Random Numbers

You can use the `runif(#, min, max)` function to generate random variables:

```
> runif(44, 50, 100)
 [1] 50.40703 74.04867 87.74888 91.64062 87.60340 99.18894 51.63613 74.95159
 [9] 78.69424 62.18938 81.59023 69.86100 62.35722 74.31295 51.53242 93.89568
[17] 74.18949 74.90164 95.80652 61.07918 77.25427 77.46069 62.96119 61.90670
[25] 64.74256 97.88017 76.18585 53.29564 64.39909 80.05913 94.13147 99.09108
[33] 98.25371 52.74621 78.27849 79.05545 82.32460 79.18416 80.67344 93.16784
[41] 78.77109 72.81770 92.80241 66.07901

> us_presidents_weights = runif(44, 50, 100)
```

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Matricies

Create matricies via the `cbind(vector**)` function:

```
> usp_matrix = cbind(us_presidents_heights, us_presidents_weights)
> usp_matrix
     us_presidents_heights us_presidents_weights
[1,]                   189              81.51802
[2,]                   170              85.93077
[3,]                   189              59.66279
[4,]                   163              74.48303
[5,]                   183              97.21214

> typeof(usp_matrix) #Returns type of matrix
[1] "double"
> class(usp_matrix) #Returns class of the object
[1] "matrix" "array"
> is.matrix(usp_matrix) #Check if usp_matrix is a matrix
[1] TRUE
> dim(usp_matrix) #Returns shape/dimensions of matrix
[1] 44  2
```
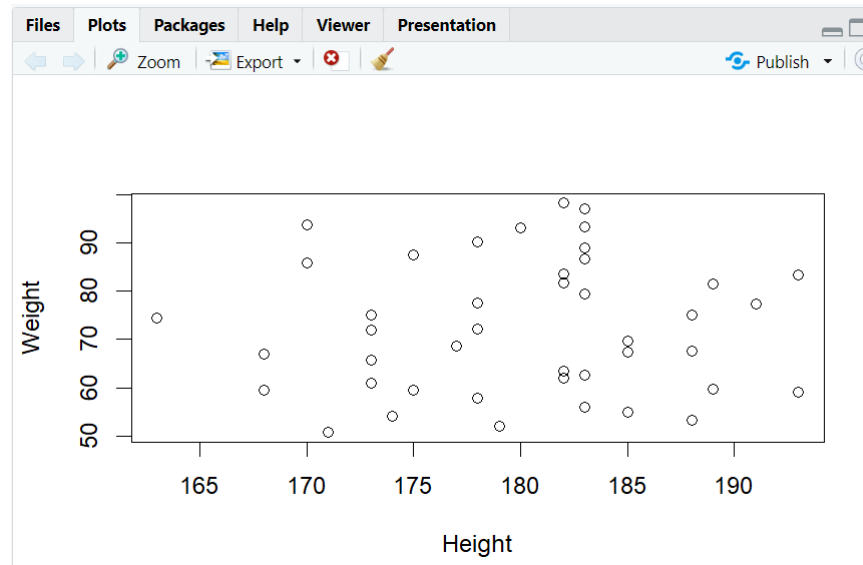
Programmierkurs 2 Data Science: R I

**Technology
Arts Sciences
TH Köln**

# Simple Plotting

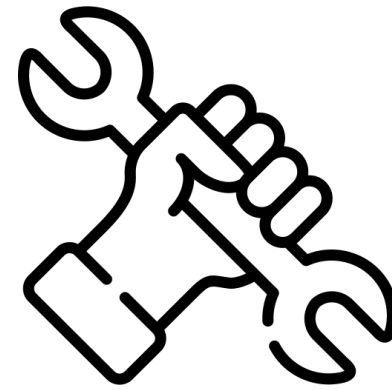Use the `plot(x,y, ylab, xlab)` function for a Scatter Plot:

```
> plot(us_presidents_heights,us_presidents_weights,ylab="Weight",xlab="Height")
```

The plot will appear in the right bottom grid under the tab „Plots":



For more advanced and attractive data visualizations, use ggplot.

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Training #2

1. Create a vector called **grades** with 100 (random) grades ranging from one to five.

2. Use the `round(vector, digits=0)` method to round the grades to integers.

3. View the frequency of grades using the `table()` function.

4. Compute the mean grade using the `length()` and `sum()` function. Do you receive the same mean grade when using the `mean()` function?

Programmierkurs 2 Data Science: R I

# Break

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# DataFrames

An extremely important data structure in R (two-dimensional table).

- Rows also called observations.

- Columns also called variables (not to be confused with the variables from before!).



variables         observations         values

R for Data Science (e2) by Wickham, Çetinkaya-Rundel, and Grolemund

Technology
Arts Sciences
TH Köln

# Reading data from .csv

Read .csv files into R using `read_csv(path)`:



```
> df_bvb_player = read_csv("https://raw.githubusercontent.com/leotraeg/FHDTM-P2DS-W
S2324/main/Data%20Science%20Projekt%20Demo/Datens%C3%A4tze/FHDTM-P2DS-WS2324-Projec
t-Demo-1.1-Data-Acquisition-Transfermarkt_BVB.csv")
Rows: 30 Columns: 8── Column specification ─────────────────────────────
───────────────
Delimiter: ","
chr (7): club_name, club_league, player_position, player_name, player_dob,...
dbl (1): player_number
ℹ Use `spec()` to retrieve the full column specification for this data.
ℹ Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

In successful .csv read; a log message tells you the

- Number of rows and columns.

- Delimiter in use.

- Column name and type specifications.

Programmierkurs 2 Data Science: R I

Technology
**Arts** Sciences
**TH Köln**

# Explore DataFrame

```
> head(df_bvb_player) #View the first rows of dataset
# A tibble: 6 × 8
  club_name    club_league player_position player_number player_name player_dob
  <chr>        <chr>       <chr>                   <dbl> <chr>       <chr>
1 Borussia D…  Bundesliga  Torwart                     1 Gregor Kob… 06.12.199…
2 Borussia D…  Bundesliga  Torwart                    35 Marcel Lot… 25.05.200…
3 Borussia D…  Bundesliga  Torwart                    33 Alexander … 13.04.199…
4 Borussia D…  Bundesliga  Torwart                    31 Silas Ostr… 19.11.200…
5 Borussia D…  Bundesliga  Abwehr                      4 Nico Schlo… 01.12.199…
6 Borussia D…  Bundesliga  Abwehr                     25 Niklas Süle 03.09.199…
# i 2 more variables: player_country <chr>, player_value <chr>
> view(df_bvb_player) #View DataFrame in an additional sub-window
```



---

Programmierkurs 2 Data Science: R I

**Technology**
**Arts Sciences**
**TH Köln**

# Explore DataFrame (cont.)

```
> nrow(df_bvb_player) #number of rows
[1] 30
> ncol(df_bvb_player) #number of columns
[1] 8
> dim(df_bvb_player) #dimension of DataFrame
[1] 30  8
> colnames(df_bvb_player) #names of columns
[1] "club_name"      "club_league"     "player_position" "player_number"
[5] "player_name"    "player_dob"      "player_country"  "player_value"

> summary(df_bvb_player) #descriptive statistics
  club_name          club_league        player_position    player_number
 Length:30          Length:30          Length:30          Min.   : 1.00
 Class :character   Class :character   Class :character   1st Qu.: 9.25
 Mode  :character   Mode  :character   Mode  :character   Median :19.50
                                                          Mean   :20.30
                                                          3rd Qu.:29.25
                                                          Max.   :47.00

  player_name         player_dob         player_country     player_value
 Length:30          Length:30          Length:30          Length:30
 Class :character   Class :character   Class :character   Class :character
 Mode  :character   Mode  :character   Mode  :character   Mode  :character
```

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Explore DataFrame (cont.)

You can access columns in a DataFrame via $ notation:

```
> df_bvb_player$player_position #access columns over $ notation
 [1] "Torwart"    "Torwart"    "Torwart"    "Torwart"    "Abwehr"
 [6] "Abwehr"     "Abwehr"     "Abwehr"     "Abwehr"     "Abwehr"
[11] "Abwehr"     "Abwehr"     "Abwehr"     "Abwehr"     "Mittelfeld"
[16] "Mittelfeld" "Mittelfeld" "Mittelfeld" "Mittelfeld" "Mittelfeld"
[21] "Mittelfeld" "Mittelfeld" "Mittelfeld" "Sturm"      "Sturm"
[26] "Sturm"      "Sturm"      "Sturm"      "Sturm"      "Sturm"
```

In the end, columns of DataFrames or nothing else than a vector (c) –
and you can apply vector functions on them such as `table()`:

```
> table(df_bvb_player$player_position)

    Abwehr Mittelfeld      Sturm    Torwart
        10          9          7          4
```

Programmierkurs 2 Data Science: R I

**Technology**
**Arts Sciences**
**TH Köln**

# Pipe %>%

**Important** tidyverse **syntax** important for Data Wrangling.

From now on, we start our coding with the DataFrame, the pipe `%>%`, and the function.

The pipe always takes what is to the left and passes it on to the function to the right:

```
> df_bvb_player %>% nrow()
[1] 30
```

In the code above, the `nrow()` function is applied to the `df_bvb_player`.

This is the same as the following line of code:

```
> nrow(df_bvb_player)
[1] 30
```

The **advantage** is that you can **connect as many functions** to the pipe **as you want**.

Programmierkurs 2 Data Science: R I

**Technology Arts Sciences TH Köln**

# DataFrame Filtering

Rows are selected using the `filter(bool_statement)` function.

The function receives one or more logical expressions as argument(s).

```
> df_bvb_player %>% filter(player_position == "Sturm")
# A tibble: 7 × 8
  club_name     club_league player_position player_number player_name player_dob
  <chr>         <chr>       <chr>                     <dbl> <chr>       <chr>
1 Borussia D… Bundesliga  Sturm                        27 Karim Adey… 18.01.200…
2 Borussia D… Bundesliga  Sturm                        43 Jamie Byno… 08.08.200…
3 Borussia D… Bundesliga  Sturm                        10 Thorgan Ha… 29.03.199…
4 Borussia D… Bundesliga  Sturm                        21 Donyell Ma… 19.01.199…
5 Borussia D… Bundesliga  Sturm                        16 Julien Dur… 05.05.200…
6 Borussia D… Bundesliga  Sturm                         9 Sébastien … 22.06.199…
7 Borussia D… Bundesliga  Sturm                        18 Youssoufa … 20.11.200…
```

Technology
Arts Sciences
TH Köln

# DataFrame Filtering (cont.)

Rows are selected using the `filter(bool_statement)` function.

The function receives one or more logical expressions as argument(s).

```
> df_bvb_player %>% filter(player_position %in% c("Abwehr","Sturm"))
# A tibble: 17 × 8
   club_name   club_league player_position player_number player_name player_dob
   <chr>       <chr>       <chr>                   <dbl> <chr>       <chr>
 1 Borussia … Bundesliga  Abwehr                      4 Nico Schlo… 01.12.199…
 2 Borussia … Bundesliga  Abwehr                     25 Niklas Süle 03.09.199…
 3 Borussia … Bundesliga  Abwehr                     15 Mats Humme… 16.12.198…
 4 Borussia … Bundesliga  Abwehr                     44 Soumaïla C… 14.10.200…
 5 Borussia … Bundesliga  Abwehr                     47 Antonios P… 10.09.199…
 6 Borussia … Bundesliga  Abwehr                      5 Ramy Bense… 16.04.199…
 7 Borussia … Bundesliga  Abwehr                     26 Julian Rye… 17.11.199…
 8 Borussia … Bundesliga  Abwehr                     17 Marius Wolf 27.05.199…
 9 Borussia … Bundesliga  Abwehr                     24 Thomas Meu… 12.09.199…
10 Borussia … Bundesliga  Abwehr                      2 Mateu More… 02.03.200…
11 Borussia … Bundesliga  Sturm                      27 Karim Adey… 18.01.200…
12 Borussia … Bundesliga  Sturm                      43 Jamie Byno… 08.08.200…
```

Technology Arts Sciences TH Köln

# DataFrame Filtering (cont.)

Rows are selected using the `filter(bool_statement)` function.

The function receives one or more logical expressions as argument(s).

```
> df_bvb_player %>% filter(player_position == "Sturm" & player_number %% 2 == 0)
# A tibble: 3 × 8
  club_name    club_league player_position player_number player_name player_dob
  <chr>        <chr>       <chr>                    <dbl> <chr>       <chr>
1 Borussia D… Bundesliga  Sturm                       10 Thorgan Ha… 29.03.199…
2 Borussia D… Bundesliga  Sturm                       16 Julien Dur… 05.05.200…
3 Borussia D… Bundesliga  Sturm                       18 Youssoufa … 20.11.200…
```

Technology
Arts Sciences
TH Köln

# DataFrame Slicing

The rows in a data frame are numbered consecutively, i.e., the rows have an **index**.

Use `slice(index or sequence)` to select rows with the internal index.

```
> df_bvb_player %>% slice(1:4)
# A tibble: 4 × 8
  club_name      club_league player_position player_number player_name player_dob
  <chr>          <chr>       <chr>                   <dbl> <chr>       <chr>
1 Borussia D…    Bundesliga  Torwart                     1 Gregor Kob… 06.12.199…
2 Borussia D…    Bundesliga  Torwart                    35 Marcel Lot… 25.05.200…
3 Borussia D…    Bundesliga  Torwart                    33 Alexander … 13.04.199…
4 Borussia D…    Bundesliga  Torwart                    31 Silas Ostr… 19.11.200…
```

Programmierkurs 2 Data Science: R I

Technology
**Arts Sciences**
**TH Köln**

# DataFrame Slicing (cont.)

> Also works on categorical attributes with alphabetical order

The functions `slice_min(column , n=1)` and `slice_max(column, n=1)` return the n rows that have the **lowest** or **highest** values in a column.

```
> df_bvb_player %>% slice_min(player_number, n=3)
# A tibble: 3 × 8
  club_name     club_league player_position player_number player_name player_dob
  <chr>         <chr>       <chr>                   <dbl> <chr>       <chr>
1 Borussia D…   Bundesliga  Torwart                     1 Gregor Kob… 06.12.199…
2 Borussia D…   Bundesliga  Abwehr                      2 Mateu More… 02.03.200…
3 Borussia D…   Bundesliga  Abwehr                      4 Nico Schlo… 01.12.199…
```

```
> df_bvb_player %>% slice_max(player_number, n=3)
# A tibble: 3 × 8
  club_name     club_league player_position player_number player_name player_dob
  <chr>         <chr>       <chr>                   <dbl> <chr>       <chr>
1 Borussia D…   Bundesliga  Abwehr                     47 Antonios P… 10.09.199…
2 Borussia D…   Bundesliga  Abwehr                     44 Soumaïla C… 14.10.200…
3 Borussia D…   Bundesliga  Sturm                      43 Jamie Byno… 08.08.200…
```

Programmierkurs 2 Data Science: R I

**Technology Arts Sciences TH Köln**

# DataFrame Selecting

To select attributes/columns/variables, you can use the function `select()`.

```
> df_bvb_player %>% select(player_name, player_position)
# A tibble: 30 × 2
   player_name          player_position
   <chr>                <chr>
 1 Gregor Kobel         Torwart
 2 Marcel Lotka         Torwart
```

Seperate multiple columns with a **comma**

```
> df_bvb_player %>% select(player_name:player_position)
# A tibble: 30 × 3
   player_name          player_number player_position
   <chr>                        <dbl> <chr>
 1 Gregor Kobel                     1 Torwart
 2 Marcel Lotka                    35 Torwart
```

Express range of columns with **colon**

```
> df_bvb_player %>% select(starts_with("player")) %>% slice(1)
# A tibble: 1 × 6
  player_position player_number player_name  player_dob      player_country
  <chr>                   <dbl> <chr>        <chr>           <chr>
1 Torwart                     1 Gregor Kobel 06.12.1997 (25) Schweiz
# i 1 more variable: player_value <chr>
```

`starts_with()` or `end_with()`

Programmierkurs 2 Data Science: R I

**Technology Arts Sciences TH Köln**

# DataFrame Mutating

We can append or change columns to data frames with `mutate()`.

It receives as arguments a new column name with the values as a vector.

```
> rep(c("Star","Rising Star", "No Star"), times=10)
 [1] "Star"        "Rising Star" "No Star"      "Star"        "Rising Star"
 [6] "No Star"     "Star"        "Rising Star" "No Star"     "Star"
[11] "Rising Star" "No Star"     "Star"        "Rising Star" "No Star"
[16] "Star"        "Rising Star" "No Star"      "Star"        "Rising Star"
[21] "No Star"     "Star"        "Rising Star" "No Star"     "Star"
[26] "Rising Star" "No Star"     "Star"        "Rising Star" "No Star"
```

```
> df_bvb_player %>% mutate(player_star_category = rep(c("Star","Rising Star", "No
Star"), times=10)) %>% select(player_name, player_star_category)
# A tibble: 30 × 2
   player_name            player_star_category
   <chr>                  <chr>
 1 Gregor Kobel           Star
 2 Marcel Lotka           Rising Star
 3 Alexander Meyer        No Star
 4 Silas Ostrzinski       Star
```

Technology
Arts Sciences
TH Köln

# DataFrame Mutating (cont.)

We can append or change columns to data frames with `mutate()`.

It receives as arguments a new column name with the values as a vector.

```
> df_bvb_player %>% mutate(number_even = ifelse(player_number %% 2 == 0,T,F)) %>%
select(player_name, player_number, number_even)
# A tibble: 30 × 3
   player_name           player_number number_even
   <chr>                         <dbl> <lgl>
 1 Gregor Kobel                      1 FALSE
 2 Marcel Lotka                     35 FALSE
 3 Alexander Meyer                  33 FALSE
 4 Silas Ostrzinski                 31 FALSE
 5 Nico Schlotterbeck                4 TRUE
```

Programmierkurs 2 Data Science: R I

**Technology**
**Arts Sciences**
**TH Köln**

# DataFrame Mutating (cont.)

We can append or change columns to data frames with `mutate()`.

It receives as arguments a new column name with the values as a vector.

```
> df_bvb_player %>% mutate(player_value_unit = ifelse(grepl("Mio", player_value),
1000000, 1000)) %>% select(player_name, player_value, player_value_unit)
# A tibble: 30 × 3
   player_name          player_value player_value_unit
   <chr>                <chr>                     <dbl>
 1 Gregor Kobel         35,00 Mio. €            1000000
 2 Marcel Lotka         1,50 Mio. €             1000000
 3 Alexander Meyer      1,00 Mio. €             1000000
 4 Silas Ostrzinski     150 Tsd. €                 1000
 5 Nico Schlotterbeck   40,00 Mio. €            1000000
```

`grepl()` searches for matches in characters

For non-binary decisions, R has the `case_when()` function.

Programmierkurs 2 Data Science: R I

**Technology**
**Arts Sciences**
**TH Köln**

# DataFrame Mutating

`mutate()` **does not change** the original **DataFrame**.

If you want, you can **overwrite** it using:

```
> #Overwrite DataFrame with assignment
> df_bvb_player = df_bvb_player %>% mutate(number_even = ifelse(player_number %% 2
== 0,T,F))
> #Overwrite DataFrame with arrow assignment
> df_bvb_player <- df_bvb_player %>% mutate(number_even = ifelse(player_number %% 2
== 0,T,F))
> #Overwrite DataFrame with double pipe assignment
> #Requires library(magrittr)
> df_bvb_player %<>% mutate(number_even = ifelse(player_number %% 2 == 0,T,F))
```

```
ifelse(cond,
value for True,
value for False)
```

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
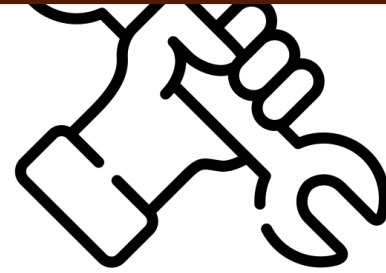TH Köln

# DataFrame Renaming

It often makes sense to **rename** columns and give them **reasonable names**.

We use the `rename(col_new = col_old)` function and overwrite using the double pipe:

```
> df_bvb_player %>% colnames()
[1] "club_name"      "club_league"      "player_position" "player_number"
[5] "player_name"    "player_dob"       "player_country"  "player_value"
[9] "number_even"


> df_bvb_player %<>% rename(player_number_even = number_even)


> df_bvb_player %>% colnames()
[1] "club_name"       "club_league"      "player_position"
[4] "player_number"   "player_name"      "player_dob"
[7] "player_country"  "player_value"     "player_number_even"
```

library(readr)

# Training #3

1. Import the following .csv dataset:
   https://raw.githubusercontent.com/leotraeg/FHDTM-P2DS-WS2324/main/Praktikum/FHDTM-P2DS-WS2324_PraktikumII_utf8.csv
   as a DataFrame called **df_dsa** in R using the readr library.

2. View the column names, dimensions, and generate a summary of df_dsa.

3. Return the frequencies of continents of the countries using table().

4. Compare whether more German students went abroad in 2015 or 2010.
   - Rename the attribute names 2015 to s_2015 and 2010 to s_2010.
   - Fill the NA values of s_2015 and s_2010 using mutate( ) and the ifelse(is.na(vector), 0, vector)) statement.
   - You should be able to use the sum( ) method to compare both years.

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Takeaways

- Know both Python and R and decide your toolkit based on your personal and other decision factors.

- R, RStudio, and tidyverse offer a broad range for data analysis.

- Similar to Python's containers, R has vectors.

- Similar to Python's DataFrame, R also has a DataFrame with filtering and slicing rows and selecting, mutating, and renaming attributes.

Programmierkurs 2 Data Science: R I

Technology
Arts Sciences
TH Köln

# Outlook

Next week we will see how to

- Deploy Functions

- Data Preprocessing

- Data Transformation

with R.

Programmierkurs 2 Data Science: R I

**Technology Arts Sciences TH Köln**

# See you again next week.

Questions?

Programmierkurs 2 Data Science: R I

**Technology Arts Sciences TH Köln**