# Introduction to R

23.11. – 25.11.2020

Christiane Hassenrück (chassenrueck@marum.de)

# Course overview

**Monday, 23.11.2020**
- General background
- Obtaining R and R studio
- R data and object types
- Data table organization, reading data into R
- Understanding errors
- Good coding practice, sustainable and collaborative coding (git and github)

**Tuesday, 24.11.2020**
- Data manipulation in base R: simple calculations, summaries, loops, functions
- Data manipulation in the R tidyverse
- Which manipulations do you need to perform with your data?

**Wednesday, 25.11.2020**
- Data visualization in R: par(), layout(), plot()
- Creating (interpolated) maps in R: *OceanView*, *marmap*, *sf*
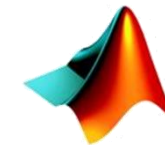- Open questions?
- Resources online and in Bremen

Course aim: To enable you to find solutions to your R problems independently!

marum

Universität Bremen

# Statistical computing software

Graphical user interface (GUI)          Command line

Commercial

Free

PAST

# Why R?

"R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS." (https://www.r-project.org/)

| Pro's | Con's |
|---|---|
| • Versatile | • Command-line |
| • Platform-independent | • Needs some getting used to… |
| • Data exploration and visualization | • Won't always tell you what to do… |
| • Hypothesis testing | |
| • Advanced graphics | |
| • Large data sets | |
| • Reproducible | |
| • Open source | |
| • Good documentation and online support | |
| … many more | |

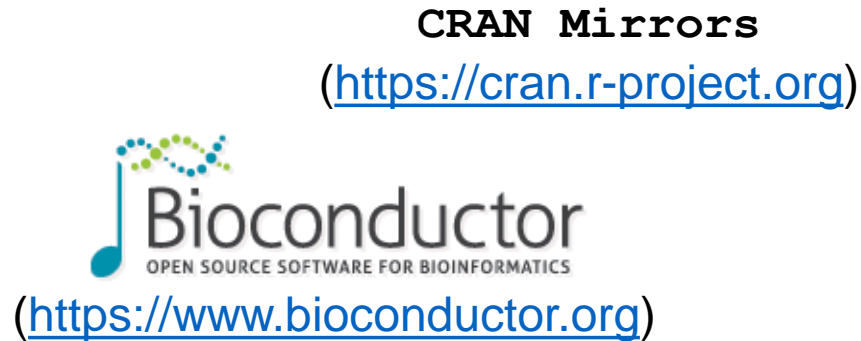**We are going to do something about that!**

# Why R?

"R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS." (https://www.r-project.org/)

R console:



Download R and R packages:

**CRAN Mirrors**
(https://cran.r-project.org)



(https://www.bioconductor.org)



(https://github.com)

More user-friendly implementation:





**R markdown**
(http://rmarkdown.rstudio.com/)

Online help:        https://stackoverflow.com/
                    https://stats.stackexchange.com/

# Why R?

# Let's get started…

Download and install R: https://cloud.r-project.org/

Download and install Rstudio: https://rstudio.com/products/rstudio/download/

# Let's get started…



**R script:**
Writing and documentation

**Console:**
Execution and (error) messages

CTRL + ENTER

**Workspace:**
R objects (data and functions)

**Various:**
File browser
Plotting panel
Help
Packages
…

# R data and object types

| | | R |
|---|---|---|
| Numerical | | |
| | Continuous (weight, temperature, length) | `numeric` |
| | Discrete (species counts) | `numeric, integer` |
| Categorical | | |
| | Ordinal (categories of increasing impact) | `ordered factor` |
| | Nominal (unweighted treatments) | `factor, character` |
| Character strings (names) | | `character` |

# R data and object types

| | **What we are used to** | **What R understands** |
|---|---|---|
| numeric | -2, 0.5, 1.452345, 50 | -2, 0.5, 1.452345, 50 |
| integer | -3, 2, 0, 6, 3, 54 | -3, 2, 0, 6, 3, 54 |
| character | sampleA, sampleB<br>exp1, exp2 | "sampleA", "sampleB"<br>"exp1", "exp2" |
| factor | exp1, exp2 | 1, 2<br>Levels: "exp1", "exp2" |
| ordered factor | exp1, exp2 | 1, 2<br>Levels: "exp1" < "exp2" |
| logical | TRUE, FALSE | TRUE, FALSE<br>1, 0 |

# R data and object types

Vector (1d)

| |
|:---:|
| "A" |
| 1.5 |
| 0.3 |
| 4 |
| -2 |
| 3.1 |
| 5 |

# R data and object types



Margin 2 (columns)

| rownames/<br>colnames | "A" | "B" | "C" |
|---|---|---|---|
| "S1" | 1.5 | -2 | 3 |
| "S2" | 0.3 | 0 | 6.6 |
| "S3" | 4 | 5 | 34 |
| "S4" | -2 | 7 | 5.2 |
| "S5" | 3.1 | 2 | -65 |
| "S6" | 5 | -89 | 0 |

Matrix (2d)

Margin 1
(rows)

# R data and object types

Margin 2 (columns)

|  rownames/<br>colnames | "A" | "B" | "C" | "D" | "E" |
|---|---|---|---|---|---|
| Data frame (2d) | | | | | |
| "S1" | 1.5 | -2 | 3 | exp1 | TRUE |
| "S2" | 0.3 | 0 | 6.6 | exp1 | TRUE |
| "S3" | 4 | 5 | 34 | exp1 | TRUE |
| "S4" | -2 | 7 | 5.2 | exp2 | TRUE |
| "S5" | 3.1 | 2 | -65 | exp2 | FALSE |
| "S6" | 5 | -89 | 0 | exp2 | FALSE |

Margin 1 (rows)

# R data and object types

List (1d): vector of R objects

|       | "A" | "B" | "C" | "D" | "E"  |
|-------|-----|-----|-----|-----|------|
| "S1"  | 1.5 | -2  | 3   | exp1| TRUE |
| "S2"  | 0.3 | 0   | 6.6 | exp1| TRUE |
| "S3"  | 4   | 5   | 34  | exp1| TRUE |
| "S4"  | -2  | 7   | 5.2 | exp2| TRUE |
| "S5"  | 3.1 | 2   | -65 | exp2| FALSE|
| "S6"  | 5   | -89 | 0   | exp2| FALSE|

data.frame

| | "A" | "B" | "C" | "D" | "E" |
|---|-----|-----|-----|-----|------|
| "S1" | 1.5 | -2 | 3 | exp1 | TRUE |
| "S2" | 0.3 | 0 | 6.6 | exp1 | TRUE |
| "S3" | 4 | 5 | 34 | exp1 | TRUE |
| "S4" | -2 | 7 | 5.2 | exp2 | TRUE |
| "S5" | 3.1 | 2 | -65 | exp2 | FALSE |
| "S6" | 5 | -89 | 0 | exp2 | FALSE |

vector

| "A" | 1.5 | 0.3 | 4 | -2 | 3.1 |
|-----|-----|-----|---|----|-----|

value

| 1.5 |
|-----|

matrix

| | "A" | "B" | "C" |
|---|-----|-----|-----|
| "S1" | 1.5 | -2 | 3 |
| "S2" | 0.3 | 0 | 6.6 |
| "S3" | 4 | 5 | 34 |
| "S4" | -2 | 7 | 5.2 |
| "S5" | 3.1 | 2 | -65 |
| "S6" | 5 | -89 | 0 |

# Let's move to R…

Exercise 1:          R data and object types

<span style="color:red">Collect unfamiliar commands!</span>

<span style="color:red">Collect your error messages!</span>

# Data table organization

- Most common input format for tabular data:
    - .txt
    - .csv
    - .tsv

- Include variable names in first row (header)
- Don't start row or column names with numbers
- Values are usually tab, space, or comma separated
- Avoid special characters and spaces in data values, variable names, and file names

|  | **Bad** | **Good** |
|---|---|---|
| Variable name | mean temperature | temperature.mean |
|  | mean-temperature | temperature_mean |
| Data value | day 1 | day1 |
|  |  | 1 (variable name: day) |

# Data table organization

The good, the bad, and the ugly…

**Merged cells**

**Hidden spaces**

**Empty cells and rows**

**Interspersed header**

**Spaces**

**Inconsistent precision**

| reef | site | seep.influence | pH | | |
|------|------|----------------|------|-------|-------|
| Illi | S1 | medium | 7.92 | 7.93 | 7.91 |
| Illi | S12 | medium | 7.94 | 7.9 | 7.99 |
| | | | | | |
| reef | site | seep.influence | SiO4 | | |
| Illi | S1 | medium | 4.47 | 4.245 | 4.956 |
| Illi | S12 | medium | 2.08 | 2.15 | 1.836 |
| | | | | | |
| reef | site | seep.influence | PO4 | | |
| Illi | S1 | medium | 0.11 | 0.107 | 0.107 |
| Illi | S12 | medium | 0.09 | 0.083 | 0.093 |

| GFF Filter | Bolinao 1 | |
|------------|-----------|---|
| | | |
| Sample Po. | Sample Name | weight mg |
| A1 | SRM1515 | |
| A2 | 219 | 129.646 |
| A3 | 177 | 128.88 |
| A4 | 210 | 125.52 |
| A5 | 202 | 131.168 |
| A6 | 91 | 134.312 |
| A7 | SRM1515 | |
| A8 | 160 | 130.936 |

| GFF Flilter | Bolinao 2 | |
|-------------|-----------|---|
| | | |
| Sample Po. | Sample Name | weight mg |
| A1 | SRM1515 | |
| A2 | 151 | 130.647 |
| A3 | 163 | 125.363 |
| A4 | 187 | 126.708 |
| A5 | 101 | 129.571 |
| A6 | 150 | 129.103 |
| A7 | SRM1515 | |
| A8 | 147 | 130.818 |

| reef | site | seep.influence | pH | SiO4 | PO4 |
|------|------|----------------|------|-------|-------|
| Illi | S1 | medium | 7.92 | 4.471 | 0.109 |
| Illi | S1 | medium | 7.93 | 4.245 | 0.107 |
| Illi | S1 | medium | 7.91 | 4.956 | 0.107 |
| Illi | S12 | medium | 7.94 | 2.076 | 0.090 |
| Illi | S12 | medium | 7.90 | 2.150 | 0.083 |
| Illi | S12 | medium | 7.99 | 1.836 | 0.093 |

# Data table organization

**Long data format:**

- One data value per line per measurement variable
- Additional comlums with contextual data (usually categories)

**Wide data format:**

- More easily readable
- Values either rearrangement of or summaries calculated from long data format

| reef | site | seep.influence | measurement | value |
|------|------|----------------|-------------|-------|
| Illi | 1 | medium | pH | 7.92 |
| Illi | 1 | medium | pH | 7.93 |
| Illi | 1 | medium | pH | 7.91 |
| Illi | 12 | medium | pH | 7.94 |
| Illi | 12 | medium | pH | 7.90 |
| Illi | 12 | medium | pH | 7.99 |
| Illi | 1 | medium | SiO4 | 4.471 |
| Illi | 1 | medium | SiO4 | 4.245 |
| Illi | 1 | medium | SiO4 | 4.956 |
| Illi | 12 | medium | SiO4 | 2.076 |
| Illi | 12 | medium | SiO4 | 2.150 |
| Illi | 12 | medium | SiO4 | 1.836 |
| Illi | 1 | medium | PO4 | 0.109 |
| Illi | 1 | medium | PO4 | 0.107 |

## Original data - rearranged

| reef | site | seep.influence | pH | SiO4 | PO4 |
|------|------|----------------|------|-------|-------|
| Illi | S1 | medium | 7.92 | 4.471 | 0.109 |
| Illi | S1 | medium | 7.93 | 4.245 | 0.107 |
| Illi | S1 | medium | 7.91 | 4.956 | 0.107 |
| Illi | S12 | medium | 7.94 | 2.076 | 0.090 |
| Illi | S12 | medium | 7.90 | 2.150 | 0.083 |
| Illi | S12 | medium | 7.99 | 1.836 | 0.093 |

## Mean values

| reef | site | seep.influence | pH | SiO4 | PO4 |
|------|------|----------------|------|-------|-------|
| Illi | S1 | medium | 7.92 | 4.539 | 0.108 |
| Illi | S12 | medium | 7.94 | 2.021 | 0.089 |

marum

Universität Bremen

# Let's move to R…

Exercise 2:          Reading data into R
                     Example data set from shallow hydrothermal vents in Papua New Guinea

Optional:            Read your own data into R (homework)


Collect your error messages!

marum                                                                    Universität Bremen

# R errors

**Syntax errors**

- When R doesn't understand you, because the command doesn't make sense…

- R returns an error message

- Majority of errors

- E.g.: Trying to calculate the mean of categorical data, typos

**Semantic errors**

- When R doesn't do what you want, although the command makes sense…

- R will not return an error message, because the command is valid

- More dangerous errors

- E.g.: Calculating percentages over columns, and not rows

**Google is your new best friend ☺**

# Sustainable and collaborative coding

Project directory organization:
└─  Input (raw) data ─────────────────────────→  Data archiving:
└─  Intermediate (temporary) files
└─  Final results (plots, tables, etc.)
└─  Scripts



(https://github.com)

Work copy ←── git clone / git pull ──── Repository

git status          git push

Untracked
        ╲ git add
Tracked
        ╲ edit
Modified
        ╲ git commit
Staged

Git plug-in for Rstudio: https://support.rstudio.com/hc/en-us/articles/200532077-Version-Control-with-Git-and-SVN

# Plot()-ing in R

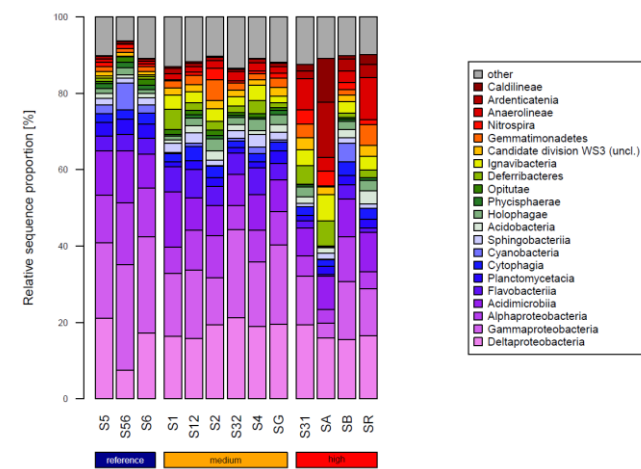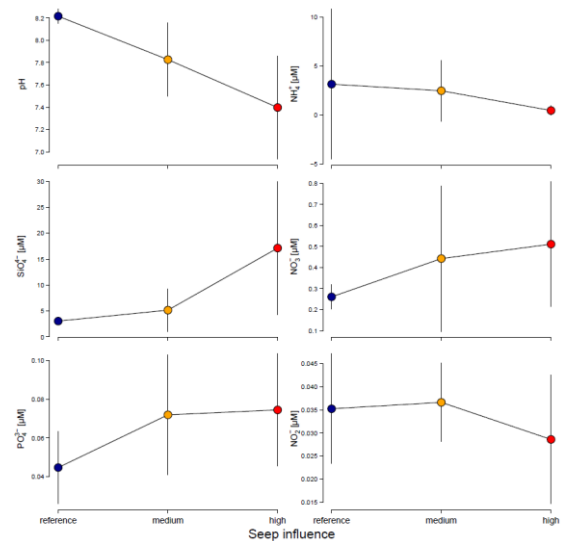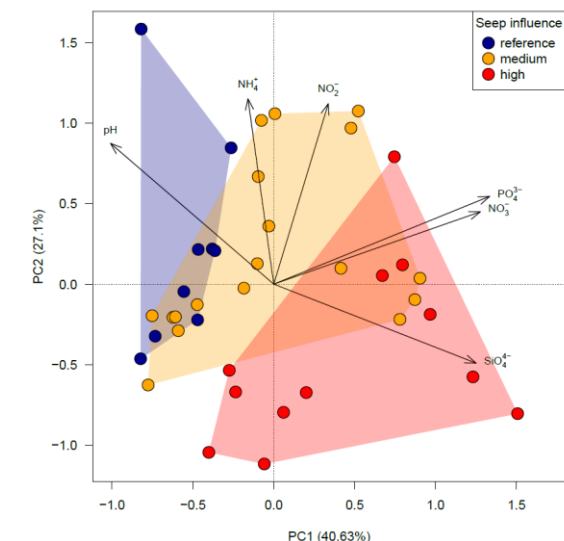Configure your plotting area:
- par()
- layout()
- plot()

Plotting elements:
- points()
- segments()
- lines()
- rect()
- polygon()
- image()
- etc.

Packages for creating maps in R:
- maps
- OceanView
- marmap
- sp, sf (more advanced spatial analysis)

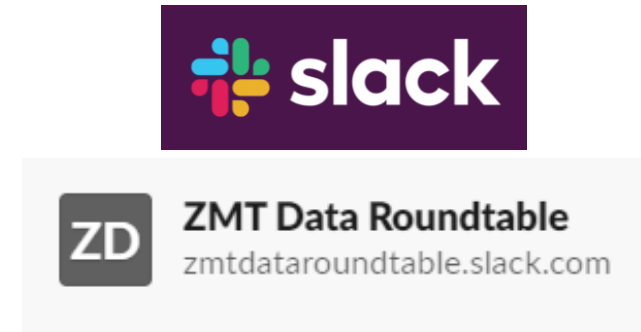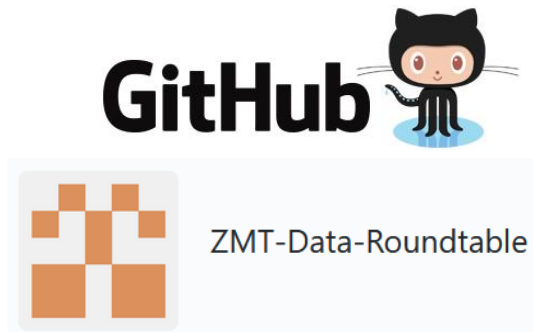Most time-consuming step when plotting:

Getting data into shape!

marum    Universität Bremen

# Example plots

# Further resources

Tidyverse tutorials: https://www.tidyverse.org/

Stackexchange: https://stats.stackexchange.com/

Stackoverflow: https://stackoverflow.com/

ZMT data roundtable: https://www.leibniz-zmt.de/de/neuigkeiten/veranstaltungen/data-round-table.html



ZMT-Data-Roundtable



ZMT Data Roundtable
zmtdataroundtable.slack.com

Contact: tobias.poprick@leibniz-zmt.de, arjun.chennu@leibniz-zmt.de

marum                                                                                     Universität Bremen