# E0410 Fundamentals of Statistics for Scientific Data Using R

by Daria Sapunova, PhD student, RECETOX

daria.sapunova@recetox.muni.cz

Bohunice, D29, room 123

# Theoretical part

# We are approaching the end of the course

| Date | Content |
|---|---|
| 25/04 | Correlation |
| 02/05 | Linear regression |
| 09/05 | Linear + multiple regression |
| 16/05 | Dean's holiday |
| 23/05 | Wrap up or linear regression if we didn't cover it up, questions |

# Final assignment

**Working on a dataset** provided **by a student or the teacher**.
The results should be presented as a **pdf file** consisting of:

| Context | Words |
|---|---|
| **a title** | - |
| **a short introduction** regarding the topic with references | approx. 200 words |
| **data description** | approx. 100 words |
| **statistical method description** | approx. 200 words |
| **results with data visualization and interpretation** of the results | approx. 200 words |
| **a short conclusion** | approx. 100 words |
| Appendix with **the R script** | - |

# Final assignment

✓ **The examination period** is from 27.05.24 to 04.07.24

✓ **Please, upload pdf files with your assignment**
  **till 01.07.24 included**

✓ **The instruction with the example** will be uploaded **18.05.-19.05.24,** prepare **your questions by 23.05**. if you will have any

✓ **The generated data** will be sent **18.05.-19.05.24**

# Who needs data?

✓ **Till 03.05.24 included fill up the table**

**Study materials -> Learning materials -> Final Assignment -> data_generate.xlsx**

| ✔ | NAME ▼ | POSTED BY | UPLOA... | RIGHTS |
|---|--------|-----------|----------|--------|
| ↑ | Final Assignment  final_assignment /3 | Sapunova, D. | Today | |
| ○ | 📁 Instructions  instructions /0 | Sapunova, D. | Today | |
| ○ | 📁 Assignments  assignments /0 | Sapunova, D. | Today | |
| ○ | 📗 data_generate.xlsx | Sapunova, D. | Today | |

# Check list

| Deadline | Action |
| --- | --- |
| 03/05 | Fill up **data_generate.xlsx** |
| 18/05-19/05 | **The instruction** with the example and **the generated data will be uploaded** by the teacher |
| 23/05 | **Questions** regarding the instruction and data |
| 01/07 | Upload the final assignment |

**If you have questions -** daria.sapunova@recetox.muni.cz

# Repetition

# Parametric and non-parametric: association between two numerical variables

**We have:** <u>two numerical continuous</u> variables



**Height**
**Weight**

**Salary**
**Productivity**

**Employment**
**Inflation**

association

association

association

# Correlation

is a measure of **the strength** of the <u>**linear association**</u>, and it **indicates the direction of the linear association.**

✓**Strength**

correlation coefficient (r/ρ) from -1 to 1

✓**Direction**

Positive or negative

✓**Direction**
Positive or negative

✓**Strength**
(r) from -1 to 1

Positive
Correlation

Negative
Correlation

No
Correlation

High Degree of Positive Correlation

High Degree of Negative Correlation

No
Correlation

# Correlation

✓**Direction**

Positive or negative

✓**Strength**

(r) from -1 to 1

| Correlation Coefficient (r) | Description (Rough Guideline ) |
|---|---|
| +1.0 | Perfect positive + association |
| +0.8 to 1.0 | Very strong + association |
| +0.6 to 0.8 | Strong + association |
| +0.4 to 0.6 | Moderate + association |
| +0.2 to 0.4 | Weak + association |
| 0.0 to +0.2 | Very weak + or no association |
| 0.0 to -0.2 | Very weak - or no association |
| -0.2 to − 0.4 | Weak - association |
| -0.4 to -0.6 | Moderate - association |
| -0.6 to -0.8 | Strong - association |
| -0.8 to -1.0 | Very strong - association |
| -1.0 | Perfect negative association |

# Correlation

✓**Direction**

Positive or negative

✓**Strength**

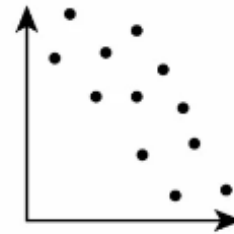(r) from -1 to 1

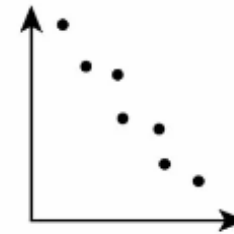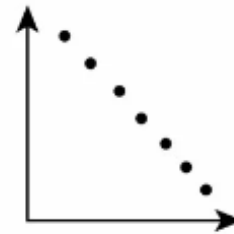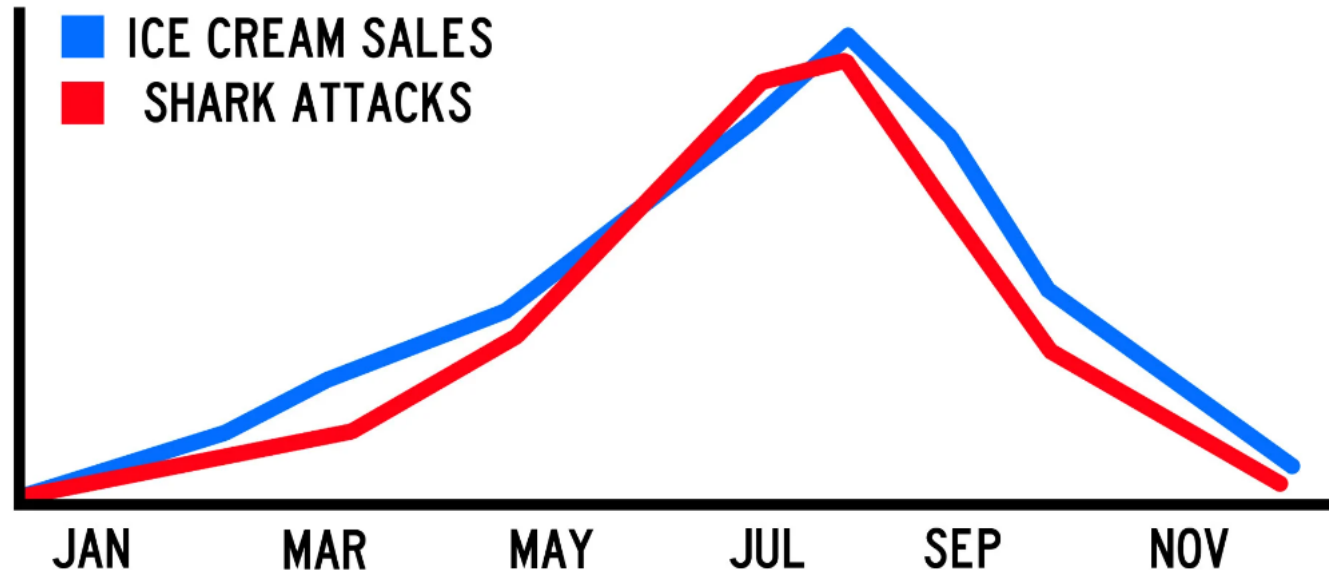| Perfect positive | Very strong positive | Moderate positive | No correlation | Moderate negative | Very strong negative | Perfect negative |
|---|---|---|---|---|---|---|



1          0.9          0.5          0          -0.5          -0.9          -1

# Correlation



## CORRELATION IS NOT CAUSATION!

■ ICE CREAM SALES
■ SHARK ATTACKS

JAN    MAR    MAY    JUL    SEP    NOV

Both ice cream sales and shark attacks increase when the weather is hot and sunny, but they are not caused by each other (they are caused by good weather, with lots of people at the beach, both eating ice cream and having a swim in the sea)

https://www.simplypsychology.org/correlation.html

# Correlation

is a measure of **the strength** of the **<u>linear association</u>**

# Parametric and non-parametric:
# association between two numerical variables

**We have:** <u>two numerical continuous</u> variables

## Parametric

✓ Pearson correlation

**! extremely sensitive to sample size,  approx.< 30**

## Nonparametric

✓ Spearman correlation
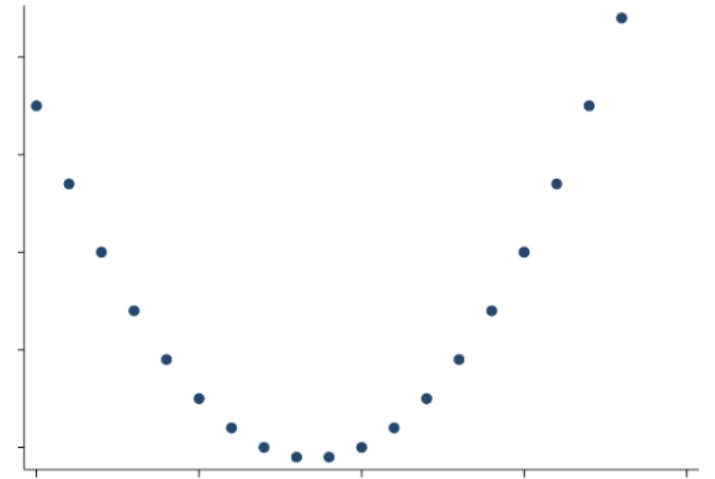
# Monotonic relationship
# between the two variables

Monotonicity means that as the value of **one variable increases**, the value of the **other variable either consistently increases or consistently decreases** (but not necessarily at a constant rate). This assumption means that the relationship between the variables **doesn't have to be strictly linear.**
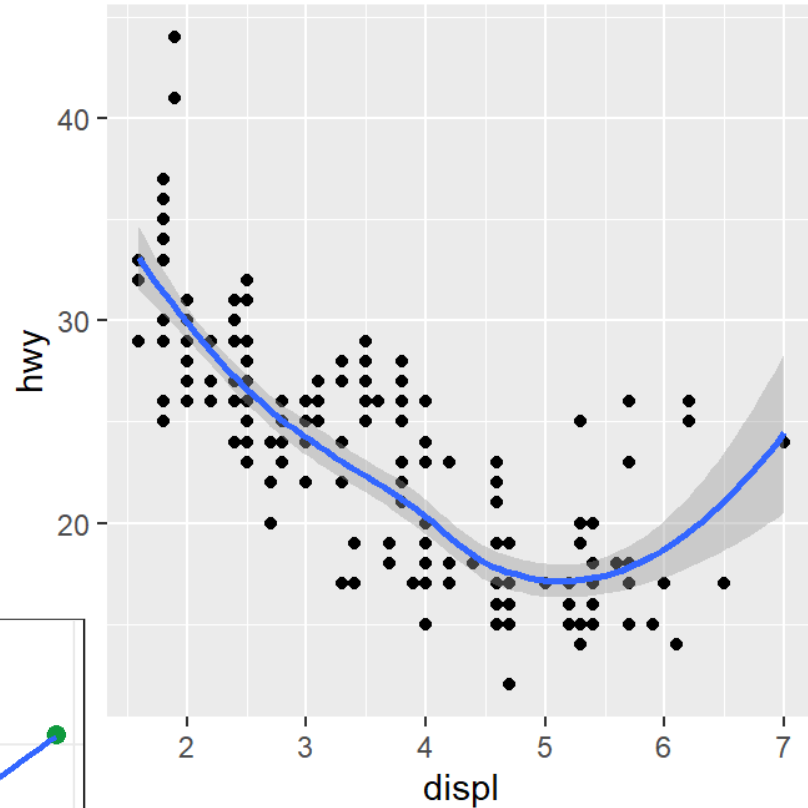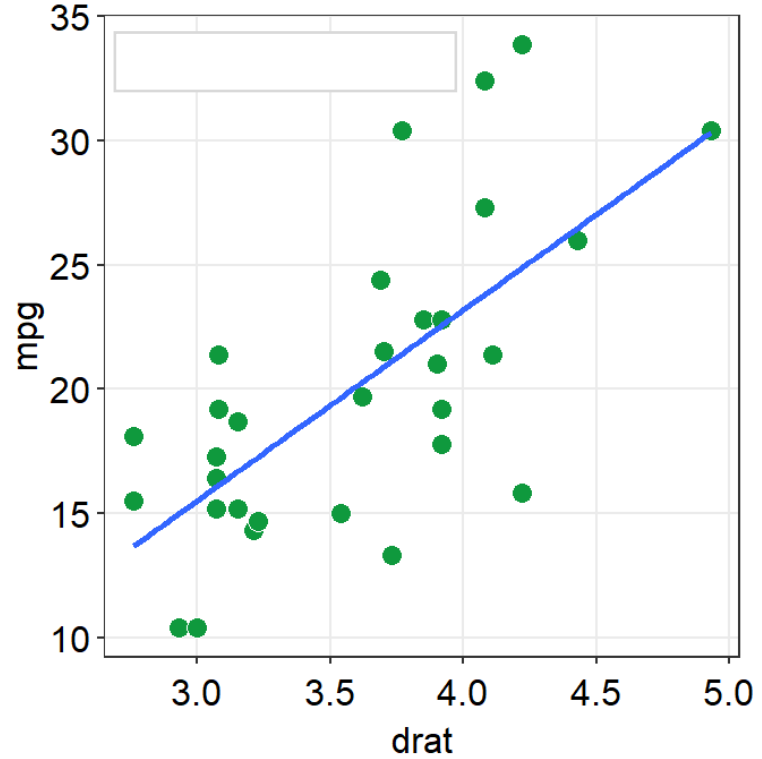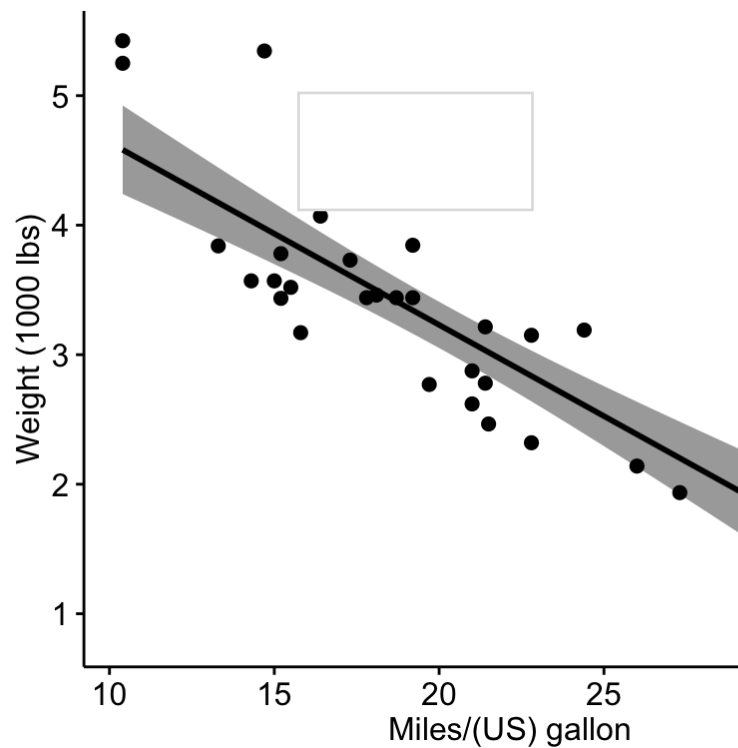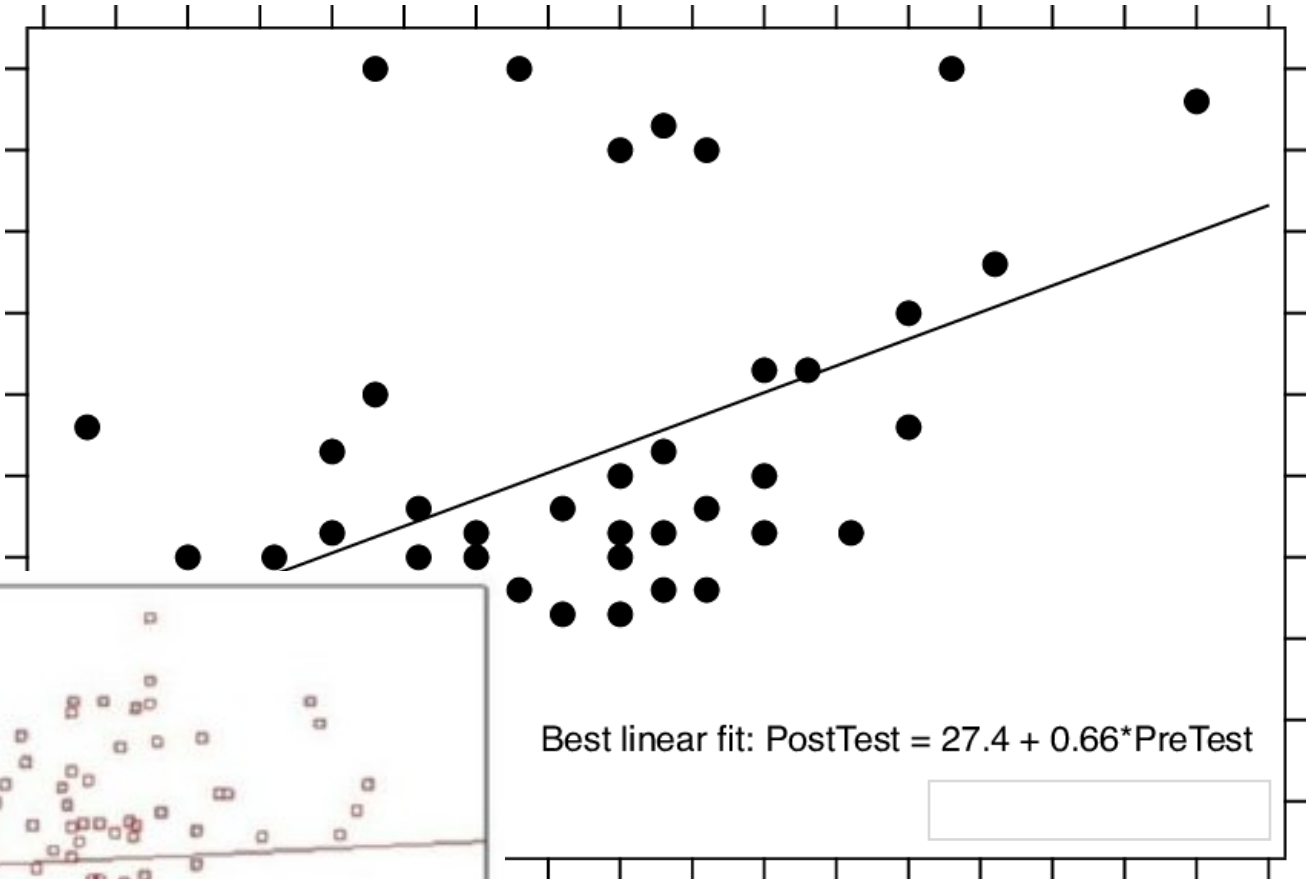


Monotonic relationship example 1

Monotonic relationship example 2

Task!

Weight (1000 lbs)

Miles/(US) gallon

mpg

drat

hwy

displ

Best linear fit: PostTest = 27.4 + 0.66*PreTest

# Practical part