# Before we start

## Some organizational details

Sarah Kaspar

Biostatistical Basics 2021

EMBL

# Course outline

EMBL

| Day | Title | Topics |
|-----|-------|--------|
| 1 | **Summarizing and visualizing data** | How to work with data frames<br>Use ggplot2 to create graphics<br>Make graphs informative |
| 2 | **Statistical distributions** | what is sampling?<br>what is a probability distribution?<br>How can we fit data to a distribution? |
| 3 | **Statistical tests** | How statistical tests work<br>Binomial test, T-test<br>Non-parametric tests |
| 4 | **Categorical data + Multiple testing** | Contingency tables<br>Test for independence<br>Measures of association<br>p-value adjustment + histogram |

**Each day:**
- lecture
- demonstration in R
- tutored exercises
- discussion of solutions

# Practical aspects

**EMBL**

**Questions:**

- allowed any time
    - unmute
    - raise hand
    - chat

**Course homepage:**

- slides
- demonstrations
- exercises

**Exercises**

- not all the functions needed are covered in the course
- take your time
- google
- ask your team mates
- ask your tutors
- solutions on the next course day
- help your collegues

# Data exploration

**Day 1**

Sarah Kaspar

Biostatistical Basics 2021

# Tools



programming language



R packages that make data

science user-friendly



Graphical user interface for R



data sets and software for analyzing

biological data

# Tidy data



"TIDY DATA is a standard way of mapping the meaning of a dataset to its structure."
—HADLEY WICKHAM

In tidy data:
- each variable forms a column
- each observation forms a row
- each cell is a single measurement

each column a variable

| id | name | color |
|----|------|-------|
| 1 | floof | gray |
| 2 | max | black |
| 3 | cat | orange |
| 4 | donut | gray |
| 5 | merlin | black |
| 6 | panda | calico |

each row an observation

Wickham, H. (2014). Tidy Data. Journal of Statistical Software 59 (10). DOI: 10.18637/jss.v059.i10

Illustrations from the Openscapes blog *Tidy Data for reproducibility, efficiency, and collaboration* by Julia Lowndes and Allison Horst

# Tidy data

**Question**: what should be the rows and columns in that table, if you

want to tidy it up?

| assessment | Billy | Suzy | Lionel | Jenny |
|---|---|---|---|---|
| quiz1 | NA | F | B | A |
| quiz2 | D | NA | C | A |
| test1 | C | NA | B | B |

| name | quiz1 | quiz2 | test1 |
|---|---|---|---|
| Billy | NA | D | C |
| Suzy | F | NA | NA |
| Lionel | B | C | B |
| Jenny | A | A | B |

# Tidy data

EMBL

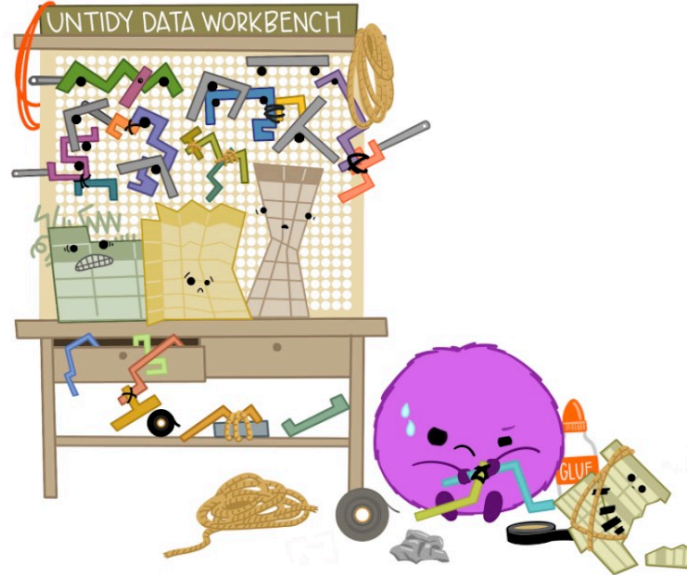| name | assessment | grade |
|------|-----------|-------|
| Billy | quiz1 | NA |
| Billy | quiz2 | D |
| Billy | test1 | C |
| Jenny | quiz1 | A |
| Jenny | quiz2 | A |
| ... | ... | ... |

- Every combinaton of name, assessment and grade is single observation.
- Every column is a variable (name, assessment, grade).
- Each cell is a single value.

# Tools for tidy data



Illustrations from the Openscapes blog *Tidy Data for reproducibility, efficiency, and collaboration* by Julia Lowndes and Allison Horst

# Motivation

EMBL

What is statistics?

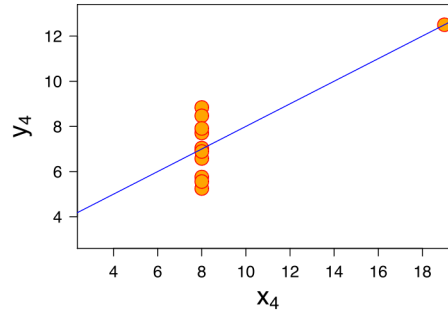A **summary statistic** "quantitatively describes or summarizes features from a collection" (Wikipedia)
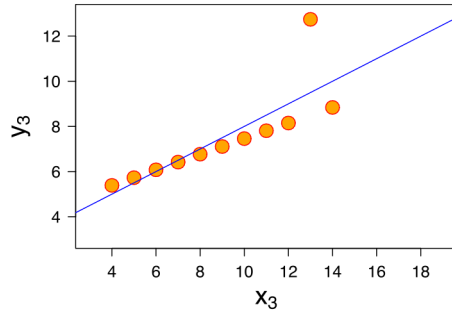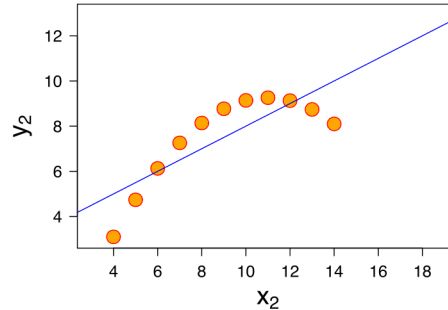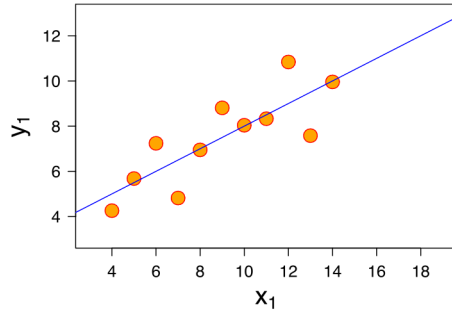
examples: mean, median, min, max,...

"**Inferential statistical analysis** infers properties of a population" (Wikipedia)

examples: hypothesis testing, t-test

# Anscombe quartett

**Why is a summary statistic not enough when exploring data?**



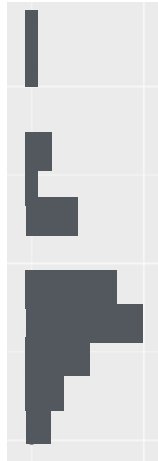All four data sets have the same mean, variance, correlation and regression line.

→ Whenever possible, plot the data points!

# Visual tools for data summary



The actual values in a distribution

How a histogram would display the values (rotated)

How a boxplot would display the values

Outliers

Whisker to farthest non-outlier point

75th percentile

50th percentile

25th percentile

1.5 x IQR

Inter-Quartile Range (IQR)

Image from: https://github.com/hadley/r4ds/blob/master/images/EDA-boxplot.pdf