# BIOSTATISTICS COURSE #1

# INTRODUCTION

SEPTEMBER 2025

# SUMMARY OF THE COURSE #1

# INTRODUCTION

01

# INTRODUCTION
## TEACHER : FRANÇOIS MACHURON

- Graduated in 2013 from **Polytech'Lille** Engineering School (GIS department)

- Engineer in Statistics and Biostatistics with 12 years of working experience

- Working experience :
  - 2013-2015 : consultant in biostatistics at **NOVARTIS** (Paris)
  - 2015-2016 : consultant in Data Science at **DECATHLON** and **COFIDIS** (Lille)
  - 2016-2019 : **Biostatistician** at CHRU Lille
  - 2020-present : **R&D Statistician** at LESAFFRE INTERNATIONAL

- Expert in Statistics / Biostatistics / Data analysis / Data mining (R & Python)

# INTRODUCTION
## BIOSTATISTICS COURSE

| Course | Date | Content |
|---|---|---|
| 1) Introduction to the course, R and descriptive statistics | 16th September 8:00 – 12:15 | History of biostatistics, definition, area of application, R & Rstudio, descriptive statistics |
| 2) Statistical modeling #1 | 23rd September 8:00 – 12:15 | Statistical tests, correlations, linear modeling, ANOVA, multiple regression |
| 3) Statistical modeling #2 | 30th September 8:00 – 12:15 | Generalized linear models, non-linear models |

# INTRODUCTION
## BIOSTATISTICS COURSE

| Course | Date | Content |
|---|---|---|
| 4) Statistical modeling #3 | 30th September 13:30 – 17:45 | PCA, MCA, FCA, clustering, discrimant analysis |
| 5) Biostatistics #1 | 7th October 8:00 – 12:15 | Levels of evidence, clinical trials, scientific publication review, effect size, power and NSN calculation |
| 6) Biostatistics #2 | 18th November 8:00 – 12:15 | Meta-analyzes, survival, longitudinal data analysis |
| 7) Exam | 25th November 8:00 – 12:15 | Exam : score /20 (with an A4 paper as help) |

# BIOSTATISTICS

02

# BIOSTATISTICS
## OVERVIEW

- Application of **statistical methods** to **biological** and **medical** data

- Many tools : **descriptive statistics**, **statistical tests**, **univariate** and **multivariate models**, **survival analysis**, NSN calculation, meta-analyzes, machine-learning...

- Many areas of application : **medicine**, **agriculture**, **biology**, **ecology**...

- More and more data collected through **clinical trials**, **R&D experiments**... => Need to apply and develop more and more biostatistical methods

- **Associated terms** : bio-informatics, data-science, epidemiology, statistics

# BIOSTATISTICS
## HISTORY

**Blaise PASCAL** (1623 – 1662) : father of probabilities and applied mathematics to science

**Pierre-Simon LAPLACE** (1749-1827) : added major contributions to probability theory and statistics

**Karl PEARSON** (1857–1936) : founder of modern statistics applied to biology and medicine

**Ronald FISHER** (1890 – 1962) : added major contributions to statistics applied to biology and medicine (ANOVA, maximum likelihood…)

# BIOSTATISTICS
## EXAMPLES OF ANALYZES

- Study of the effect of an **experimental treatment on a disease** (example : Pfizer-BioNTech COVID-19 vaccine effect on infection, hospitalization and mortality)

- **Comparison** of baseline parameters between groups with statistical tests

- **Meta-analysis** : systematic method that uses statistical techniques for combining results from different studies to obtain a quantitative estimate of the overall effect of a particular intervention or variable on a defined outcome

- **Logistic regression** : construction of a predictive model on smoking status from socio-demographic data of a cohort

- **Survival analysis** : comparison of effect of a treatment on mortality

# EXAMPLE OF A CLINICAL TRIAL

03

# EXAMPLE OF A CLINICAL TRIAL
## PLEASE READ THE ARTICLE (30 MINUTES)

# EXAMPLE OF A CLINICAL TRIAL
## OVERVIEW

- **Objective** : Prevention of Malaria in HIV-uninfected Pregnant **Women** and **Infants** in **Uganda**

- **Type of study** : Double blinded **randomized** controlled trial (link to the trial : https://clinicaltrials.gov/study/NCT02793622)

- **Intervention** : Monthly **Sulfadoxine-pyrimethamine** (SP) or **Dihydroartemisinin-piperaquine** (DP)

- **Hypothesis :** DP will significantly reduce the **burden of malaria** in pregnancy and infancy compared to SP

# EXAMPLE OF A CLINICAL TRIAL
## OVERVIEW

- **782 pregnant women** included (391 in each group) (12 to 20 weeks of pregnancy)

- **Two treatments** taken monthly during pregnancy : Sulfadoxine-pyrimethamine (SP) vs Dihydroartemisinin-piperaquine (DP)

- Study conducted between **September 2016 and May 2017.**

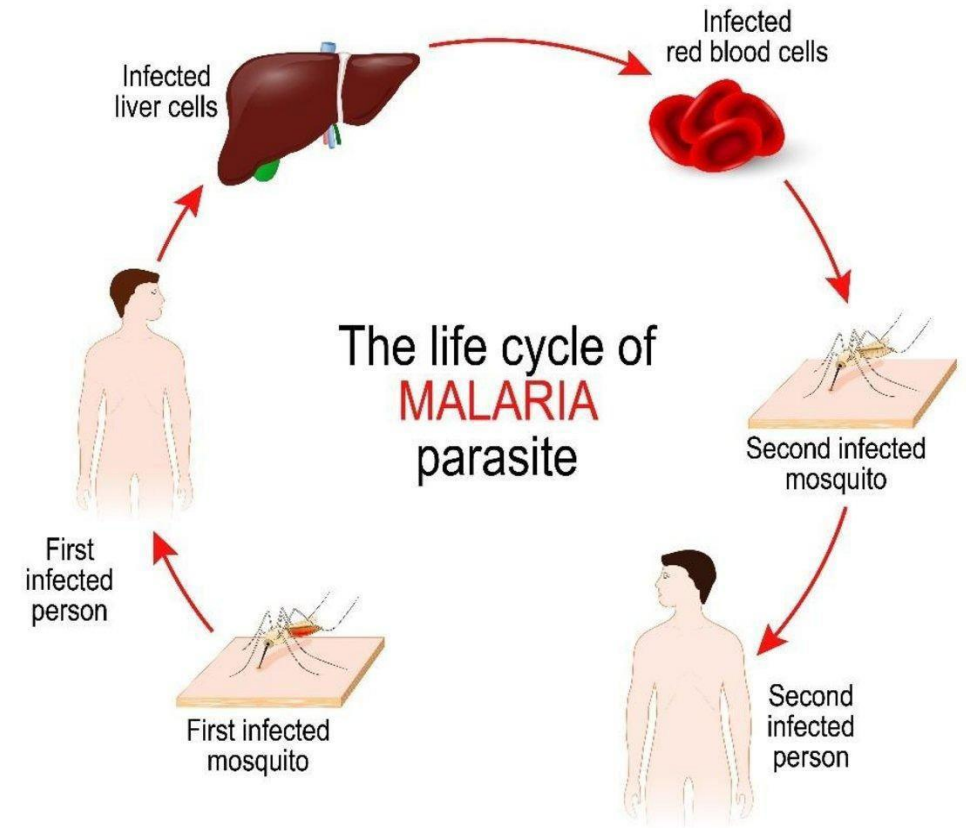- Results published in the **Lancet on 22nd March 2019**



**Figure 1: Trial profile**
DP=dihydroartemisinin–piperaquine. SP=sulfadoxine–pyrimethamine.

# EXAMPLE OF A CLINICAL TRIAL
## MALARIA DISEASE

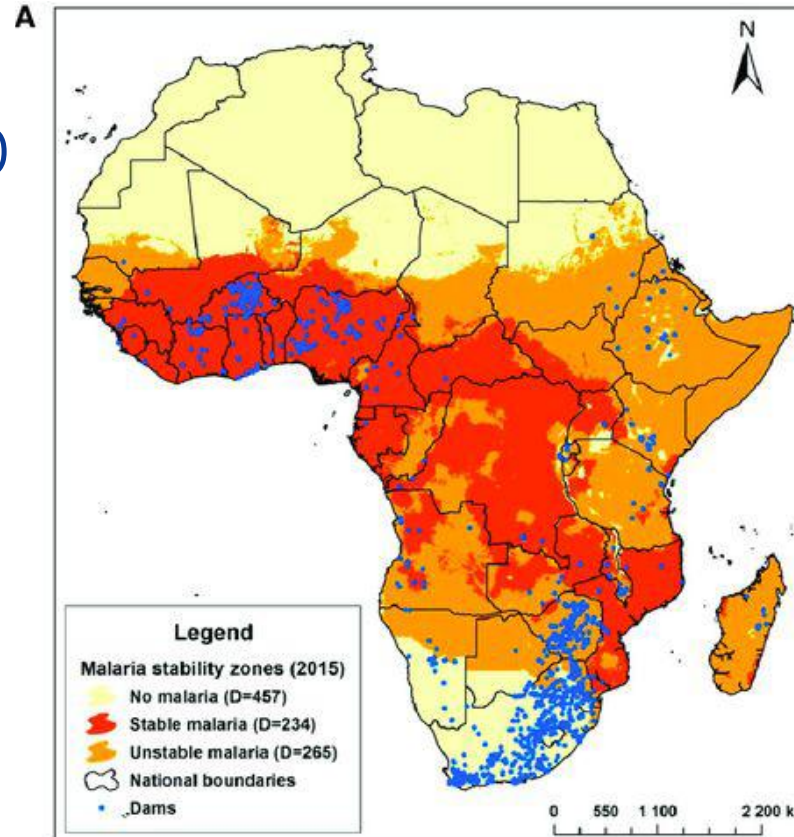- Caused by single-cell **parasites** that at spread by **mosquitoes**

- Damages especially **liver** and **red blood cells**

- Causes waves of **fevers** every few days, extreme **anemia** and **blood clots**

- May lead to "*complicated malaria*", which is described as **sepsis** which can lead to **death**

# EXAMPLE OF A CLINICAL TRIAL
## MALARIA EPIDEMIOLOGY

- **Hundreds of millions** of cases worldwide, **660,000** deaths each year

- **81% of cases and 91% of deaths are in Africa**

- **Pregnant women** and children under 5 most vulnerable



- **Prevention in Africa** : Regular or insecticide treated **bednets** and preventative therapy with **Sulfadoxine-pyrimethamine** (SP)

# EXAMPLE OF A CLINICAL TRIAL
## MALARIA DURING PREGNANCY

**Maternal outcomes**

- Maternal anaemia
- Cerebral malaria
- Severe malaria
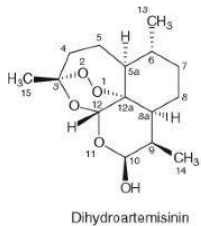- Maternal mortality
- Recurrent or new plasmodium infections

**Child outcomes**

- Abortion
- Stillbirth
- Preterm delivery
- Low birth weight
- Neonatal mortality
- Congenital malaria
- Infant mortality
- Anaemia
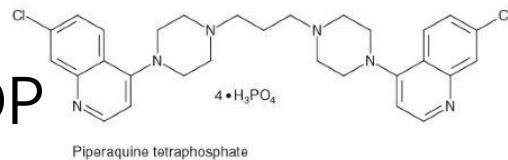- Poor developmental / behavioural outcome

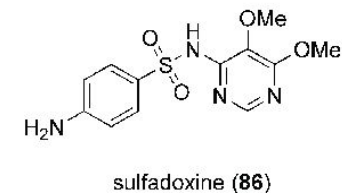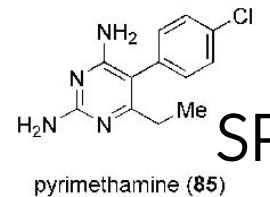# EXAMPLE OF A CLINICAL TRIAL
## MOTIVATIONS

- **Antifolate resistance** of malaria parasite (65% in Eastern Africa) : Sulfadoxine-pyrimethamine (SP) does not work anymore

- **Artemisinin effective** alternative treatment for malaria and is safe for treatment of malaria in pregnant women

- Dihydroartemisinin-piperaquine (DP) together with bednets associated with **a lower incidence of malaria infection during pregnancy and after delivery**

- **Serious adverse events lower** with DP than SP



DP

versus

SP

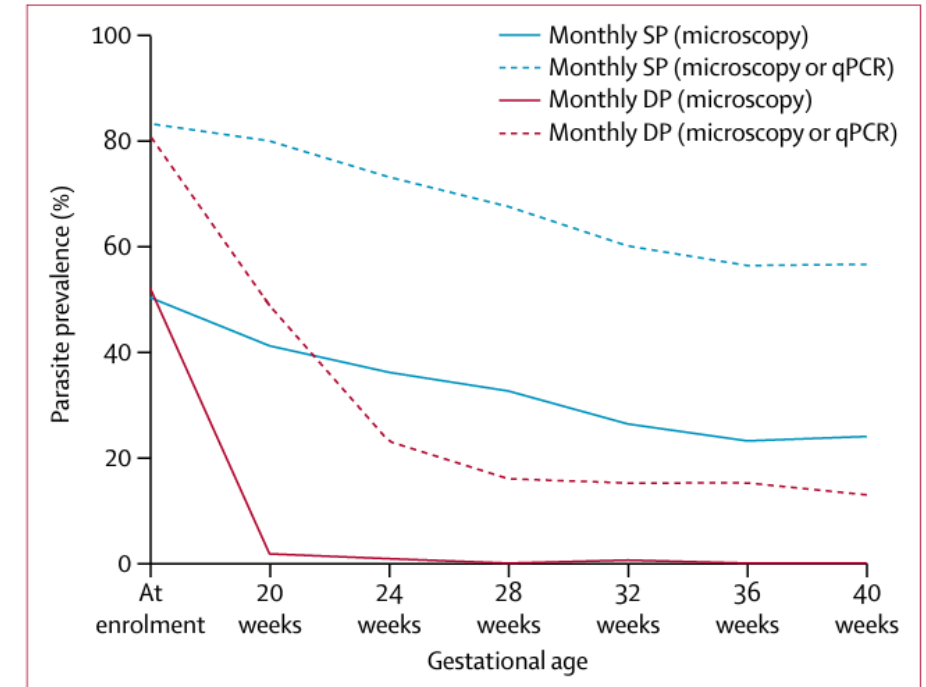# EXAMPLE OF A CLINICAL TRIAL
## RESULTS – BASELINE CHARACTERISTICS

| | Monthly sulfadoxine–pyrimethamine group (n=391) | Monthly dihydroartemisinin–piperaquine group (n=391) |
|---|---|---|
| Age, years | 23 (19–27) | 23 (19–27) |
| Gestational age, weeks | 15·4 (13·3–17·6) | 15·0 (13·4–17·1) |
| Gestational age category (weeks) | | |
| 12–16 | 234 (60%) | 242 (62%) |
| >16–20 | 157 (40%) | 149 (38%) |
| Gravidity | | |
| 1 | 102 (26%) | 93 (24%) |
| 2 | 85 (22%) | 105 (27%) |
| ≥3 | 204 (52%) | 193 (49%) |
| ITN coverage | | |
| Ownership of ITN | 61 (16%) | 50 (13%) |
| Reported sleeping under an ITN the previous night | 48 (12%) | 44 (11%) |
| Haemoglobin concentration (g/dL) | 11·5 (1·3) | 11·4 (1·2) |
| Detection of malaria parasites by microscopy | 197 (50%) | 204 (52%) |
| Detection of malaria parasites by microscopy or qPCR | 326 (83%) | 317 (81%) |

Data are median (IQR) or n (%). ITN=insecticide-treated net. qPCR=quantitative PCR.

*Table 1:* Baseline characteristics

# EXAMPLE OF A CLINICAL TRIAL
## RESULTS - OUTCOMES

| | Monthly sulfadoxine–pyrimethamine group* (n=338) | Monthly dihydroartemisinin–piperaquine group (n=349) | Protective efficacy† (95% CI) | p value |
|---|---|---|---|---|
| **Outcomes assessed at delivery** | | | | |
| Adverse birth outcomes among livebirths | | | | |
| Composite (primary efficacy outcome) | 60/329 (18%) | 54/337 (16%) | 12% (–23 to 37) | 0·45 |
| Low birthweight | 29/329 (9%) | 24/337 (7%) | 19% (–36 to 52) | 0·42 |
| Preterm birth | 24/329 (7%) | 16/337 (5%) | 35% (–20 to 65) | 0·17 |
| Small for gestational age | 41/329 (13%) | 39/337 (12%) | 7% (–40 to 38) | 0·72 |
| Fetal or neonatal death | | | | |
| Spontaneous abortion | 4/338 (1%) | 10/349 (3%) | –142% (–665 to 23) | 0·13 |
| Stillbirth | 5/334 (2%) | 2/339 (1%) | 61% (–102 to 92) | 0·26 |
| Neonatal death | 6/329 (2%) | 4/337 (1%) | 35% (–129 to 81) | 0·50 |
| Composite | 15/338 (4%) | 16/349 (5%) | –3% (–106 to 48) | 0·93 |
| Measures of infection with malaria parasites | | | | |
| Maternal blood positive for malaria parasites by microscopy | 28/336 (8%) | 1/342 (<1%) | 96% (74 to 99) | 0·0010 |
| Placental blood positive for malaria parasites by microscopy | 28/320 (9%) | 1/333 (<1%) | 97% (75 to 99) | 0·0009 |
| Placental blood positive for malaria parasites by LAMP | 70/312 (22%) | 7/328 (2%) | 90% (80 to 96) | <0·0001 |
| Placental tissue positive for malaria parasites or pigment | 197/322 (61%) | 94/331 (28%) | 54% (44 to 62) | <0·0001 |
| **Incidence measures during pregnancy‡** | | | | |
| Symptomatic malaria | 75 (0·52)§ | 3 (0·02)§ | 96% (88 to 99) | <0·0001 |
| **Prevalence measures during pregnancy‡** | | | | |
| Detection of malaria parasites by microscopy | 519/1687 (31%) | 9/1757 (1%) | 98% (96 to 99) | <0·0001 |
| Detection of malaria parasites by microscopy or qPCR | 1105/1676 (66%) | 369/1746 (21%) | 68% (64 to 71) | <0·0001 |
| Anaemia (haemoglobin <10 g/dL) | 171/870 (20%) | 89/904 (10%) | 50% (32 to 73) | <0·0001 |

LAMP=loop-mediated isothermal amplification. qPCR=quantitative PCR. *Reference group. †Protective efficacy=1–incidence rate ratio or 1–prevalence ratio. ‡Assessed at the time of routine visits following administration of first dose of study drugs. §Number of events (incidence per person-year at risk).

*Table 2:* Efficacy outcomes



*Figure 2:* Parasite prevalence during pregnancy according to week of gestation
Parasite prevalence was assessed by microscopy and qPCR. Data at 20 weeks gestational age only includes the subset of women who received their first dose of study drugs at 16 weeks gestational age. DP=dihydroartemisinin–piperaquine. qPCR=quantitative PCR. SP=sulfadoxine–pyrimethamine.

R & RSTUDIO

04

# R & RSTUDIO
## BENCHMARK OF STATISTICAL SOFTWARES

### Microsoft Excel
- Spreadsheet software
- Simple mathematical operations
- Graphical UI

### Statistical Analysis Software
- "Official" software for clinical trials (FDA...)
- Simple and complex statistics
- Graphical UI & coding language

### R
- Collaborative language (lot of libraries available)
- **Particularly suitable in biology and medical statistics**
- Graphical UI (Rstudio) & coding language

### IBM SPSS
- Widely used
- Simple statistics and visualization
- Graphical UI

### MATLAB
- Complex mathematical processes
- Suitable for engineering studies
- Own coding language

### PYTHON
- General-purpose programming language
- Most used for machine-learning and AI
- Own coding language

# R & RSTUDIO
## MOTIVATIONS TO USE

- **Easy-to-learn** and use language (somehow like Python language)

- Strong academical and industrial **community**

- **Thousands of packages** validated, maintained and updated packages specialized in data analysis, machine-learning, AI...particularly for biological and medical data : available on the **CRAN** website (https://cran.rstudio.com/)

- **Well documented** language and packages

- Useful for **data visualization** (ggplot2, plotly, Rmarkdown and Shiny packages)

- Fully integrated by many **LLMs** (Copilot, ChatGPT, Gemini, Claude Code...)

# R & RSTUDIO
## MOTIVATIONS TO USE R Studio®

- *"Official"* R IDE (Integrated Development Environment)

- Free and maintained (current version 2024.12.0+467 released on 16th Dec 2024)

- Clear and customizable interface

- Easy to connect to **databases** and **cloud providers** (AWS, Azure, Google...)

- Able to run **Python code**

- Useful for **web-app development** (Shiny package)

- Download here on the site of POSIT

# R & RSTUDIO
## INTERFACE OF R Studio®



Coding panel

Environment components

Console & Terminal

Files
Plots
Packages
Help

# R & RSTUDIO



Live demo

# R & RSTUDIO

Time to play !
(15 minutes)

# USEFUL
# PACKAGES

# 05

# USEFUL PACKAGES
## PACKAGE DPLYR

- Powerful package for **data transformation, manipulation & summarization**

- Allow to use **pipes** (%>%) : efficient chain of operations on data

- Main functionalities :
  - ➢ Function **mutate** : new variables creation
  - ➢ Function **select** : keep or drop variables from a dataset
  - ➢ Function **filter** : select specific rows in a dataset
  - ➢ Function **arrange** : sort a dataset on one or many variables
  - ➢ Function **summarise** : statistics calculation
  - ➢ Function **group_by** : group rows (useful before summarise function)

- How to install in R : **install.packages("dplyr")**

# USEFUL PACKAGES
## PACKAGE DPLYR



Live demo

dplyr_cheatsheet.pdf

# Time to play !
# (15 minutes)

# USEFUL PACKAGES
## PACKAGE STRINGR

- Powerful package for **string processing**

- Main functionalities :
  - ➢ Function **str_detect** : pattern detection in a string (return True or False)
  - ➢ Function **str_locate** : locates the first position of a pattern and returns a numeric matrix with columns start and end
  - ➢ Function **str_extract** : extracts text corresponding to the first match, returning a character vector
  - ➢ Function **str_replace** : replaces the first matched pattern and returns a character vector

- How to install in R : install.packages("stringr")
- Regular expressions : regular expressions with stringR

# USEFUL PACKAGES
## PACKAGE STRINGR


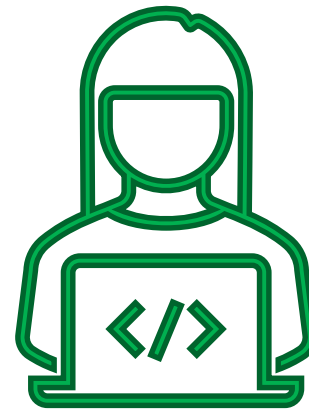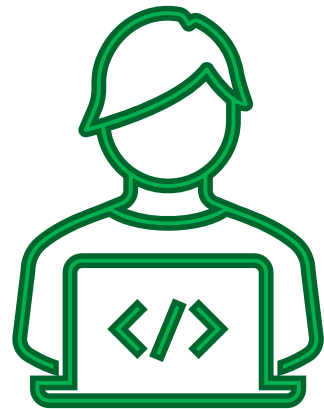
Live demo

stringr_cheatsheet.pdf

# USEFUL PACKAGES
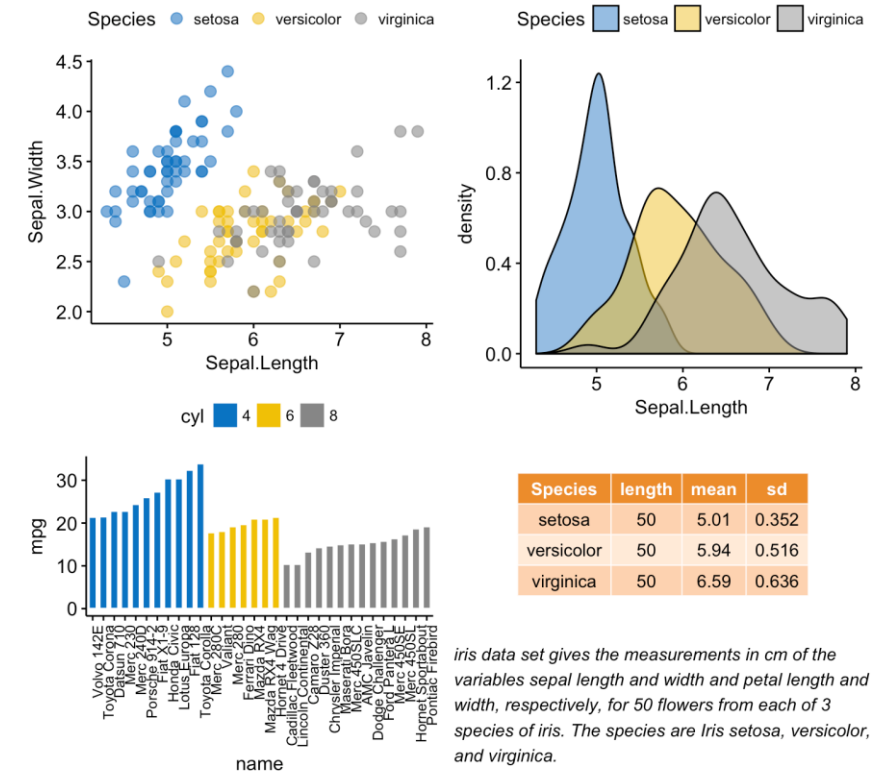## PACKAGE STRINGR

Time to play !
(15 minutes)

# USEFUL PACKAGES
## PACKAGE LUBRIDATE

- Powerful package for **dates and times processing**

- **Date-time format in R** : 2025-01-22 10:30:46 (stored in R in as the number of seconds since 1970-01-01 00:00:00 UTC)

- Main functionalities :
  - ➤ Function **date** : extract the date component of a date-time variable
  - ➤ Functions **year, month or day** : extract a specific component of date of a date-time variable
  - ➤ Function **hour, minute or second** : extract a specific component of time of a date-time variable

- How to install in R : **install.packages("lubridate")**

# USEFUL PACKAGES
## PACKAGE LUBRIDATE



## Live demo

lubridate_cheatsheet.pdf

# USEFUL PACKAGES
## PACKAGE LUBRIDATE

Time to play !
(15 minutes)

# USEFUL PACKAGES
## PACKAGE GGPLOT2



- Powerful package for **data visualization** with highly customizable **static plots**
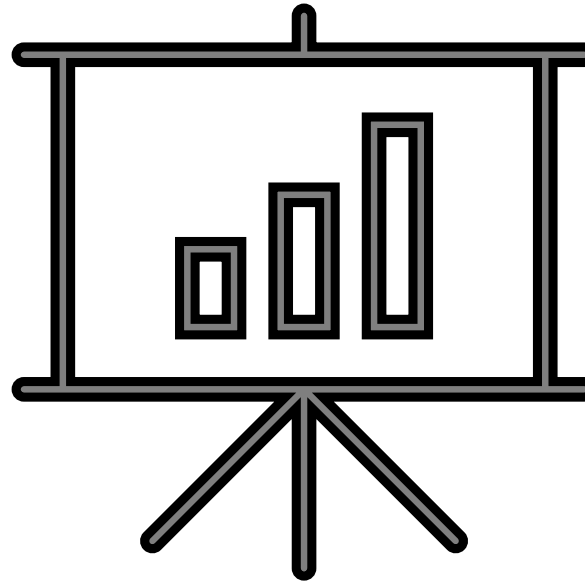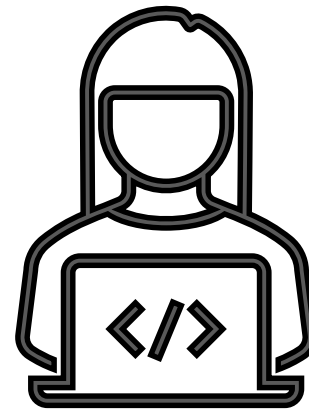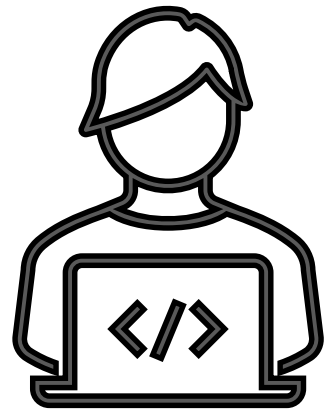
- Uses a **pipe-wise coding grammar** :

      ggplot(data=dataset) +
         element_1()+
         element_2()+
         ...

- **Elements are added** on the same plot **with layers** (dozen of graphical elements available)



- How to install in R : install.packages("ggplot2")

# Live demo

ggplot2_cheatsheet.pdf

# USEFUL PACKAGES
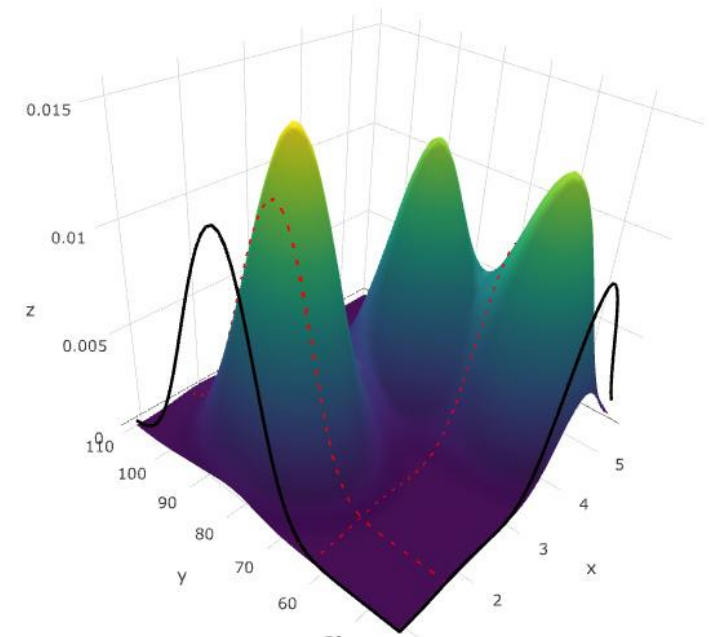## PACKAGE GGPLOT2

Time to play !
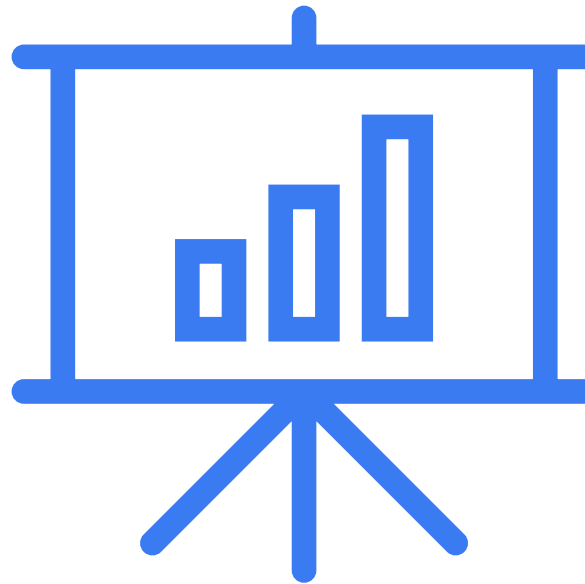(30 minutes)

# USEFUL PACKAGES
## PACKAGE PLOTLY



- Powerful package for **data visualization** with highly customizable **interactive plots**

- Available also in Python with a similar syntax

- Easily integrated in websites, web-apps or Rmarkdown HTML documents

- Useful with dplyr pipes (%>%) : elements are added on the same plot **with layers** (dozen of graphical elements available)

- How to install in R : install.packages("plotly")
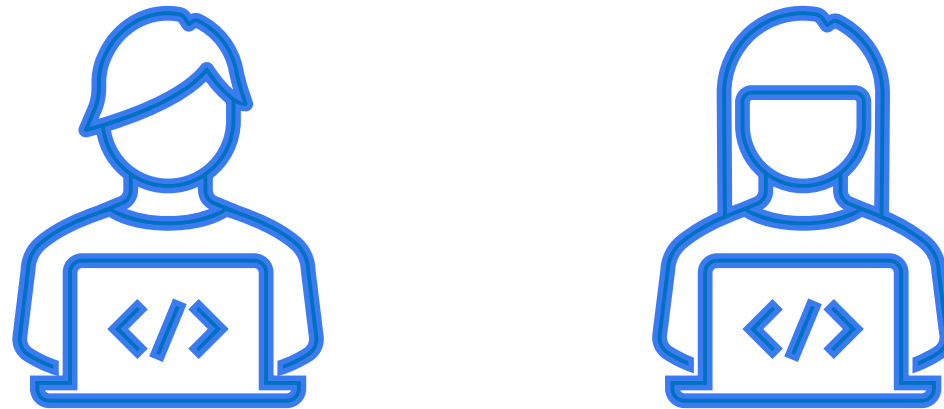
# USEFUL PACKAGES
## PACKAGE PLOTLY



Live demo

plotly_cheatsheet.pdf

# USEFUL PACKAGES
## PACKAGE PLOTLY

**Time to play !**
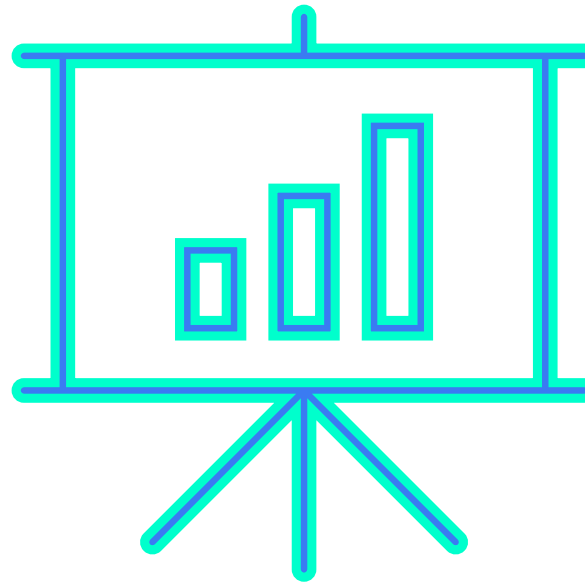**(20 minutes)**

# USEFUL PACKAGES
## PACKAGE RMARKDOWN

- Powerful package for **highly customizable and interactive documents**

- Can use multiple languages (R, Python, SQL)

- Customizable with HTML, CSS languages

- Can integrate interactive plots (produced with Plotly package)

- Output : Microsoft Word, Powerpoint documents, HTML file, PDF file

- How to install in R : install.packages("rmarkdown")
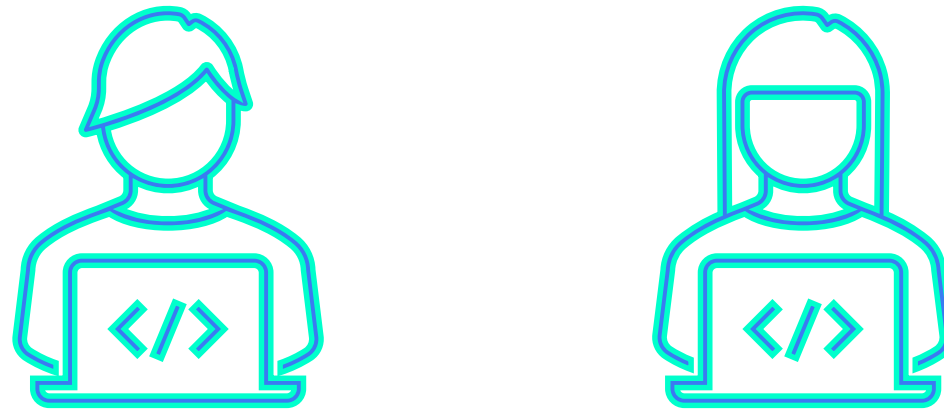
# USEFUL PACKAGES
## PACKAGE RMARKDOWN
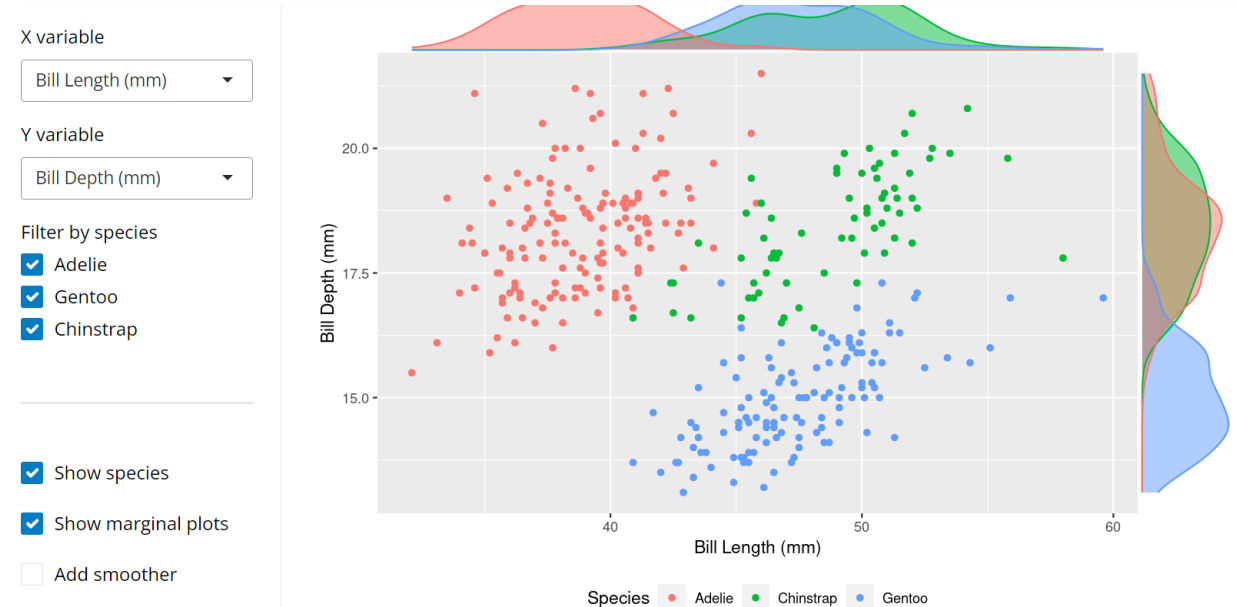
Live demo

# USEFUL PACKAGES
## PACKAGE RMARKDOWN

Time to play !
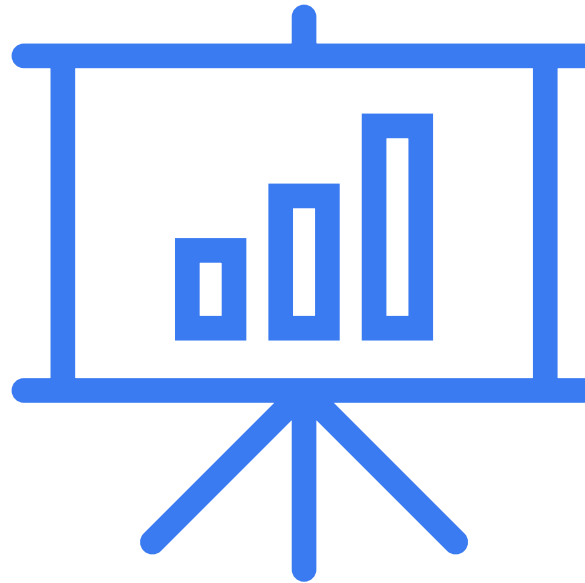(20 minutes)

# USEFUL PACKAGES
## PACKAGE SHINY

- Powerful package for **web-app development**

- Available also in Python with a similar syntax

- Many input elements (buttons, sliders, selectors, panels...)

- Customizable with HTML, Javascript languages

- How to install in R : install.packages("shiny")

# USEFUL PACKAGES
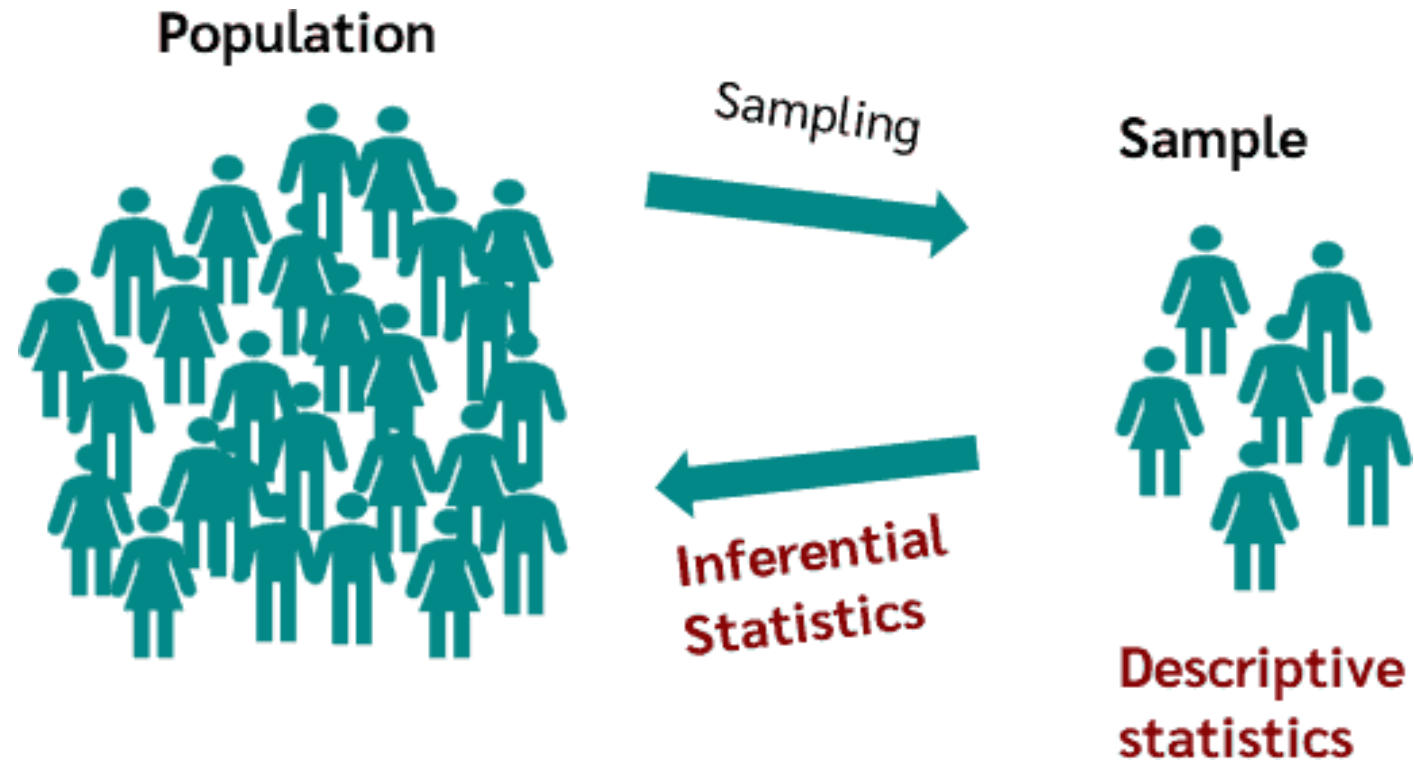## PACKAGE SHINY



# Live demo

# DESCRIPTIVE STATISTICS

06

# DESCRIPTIVE STATISTICS
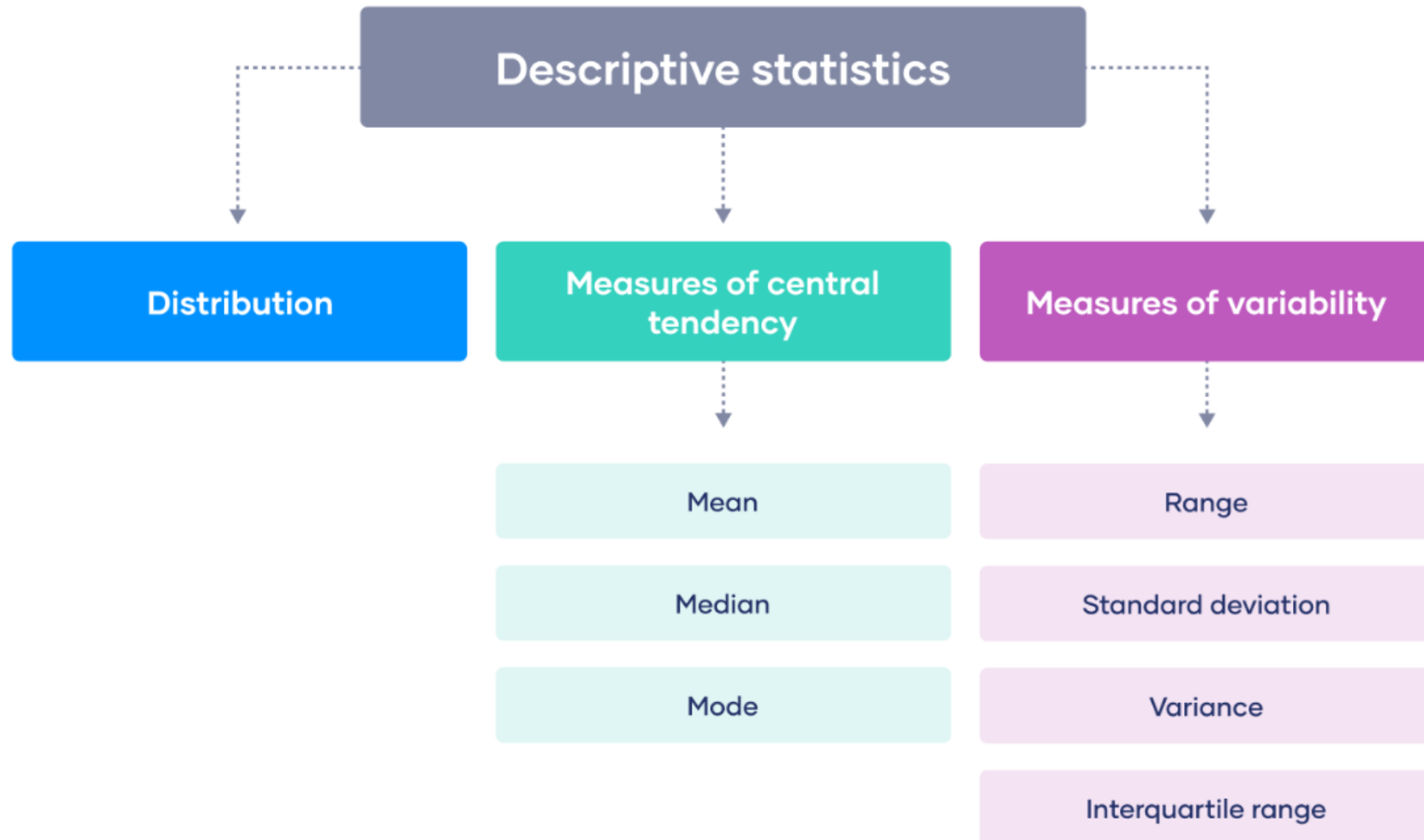## DESCRIPTIVE VS INFERENTIAL STATISTICS

# DESCRIPTIVE STATISTICS
## DESCRIPTIVE VS INFERENTIAL STATISTICS

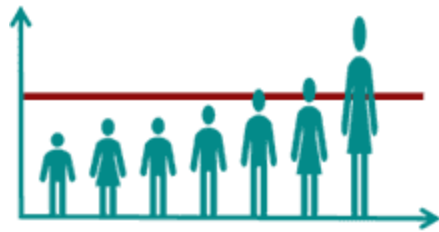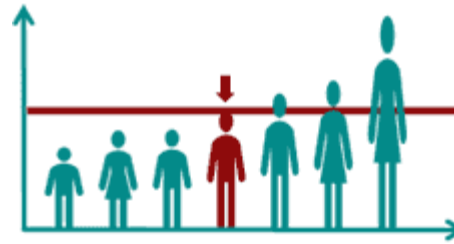| DESCRIPTIVE | INFERENTIAL |
|---|---|
| It is the analysis of data that helps to describe, show and summarize data under study | It is the analysis of random sample of data taken from a population to describe and make inference about the population |
| Organize, analyze and present data in a meaningful way | Compares, test and predicts data |
| It is used to describe a situation | It is used to explain the chance of occurrence of an event |
| It explain already known data and limited to a sample or population having small size | It attempts to reach the conclusion about the population |
| Types: Measure of central tendency & Measure of variability | Types: Estimation of parameters & Testing of hypothesis |
| Results are shown with help of charts, graphs, tables etc. | Results are shown with help of probability scores |

# DESCRIPTIVE STATISTICS
## OVERVIEW

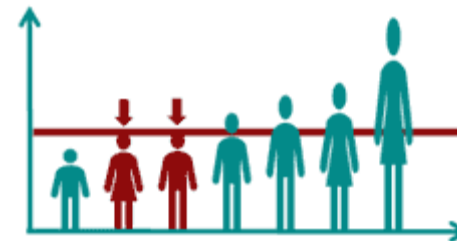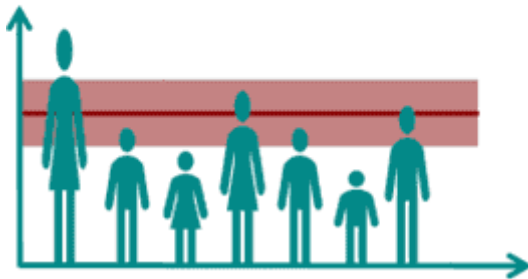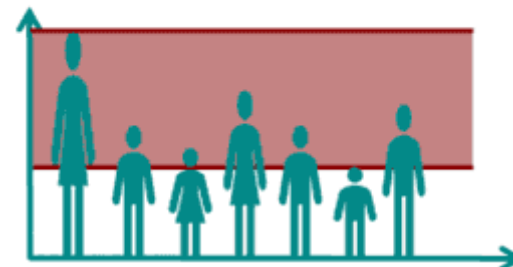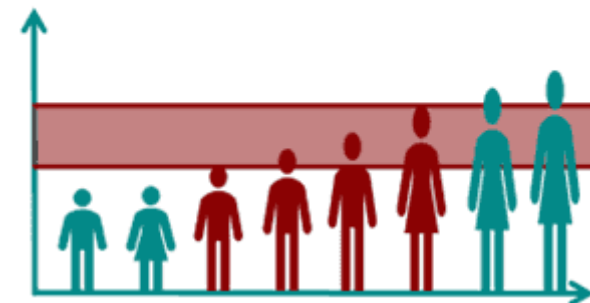# DESCRIPTIVE STATISTICS
## OVERVIEW



Mean

Median

Mode

Standard Deviation

Range

Interquartile Range

# DESCRIPTIVE STATISTICS
## FUNCTIONS FOR CONTINUOUS VARIABLES

The functions to use in R are included in the *stats* package :

- Arithmetic mean $\hat{x}$ : *mean* function
  Parameters :    x = list of continuous values
  
  $$\hat{x} = \frac{1}{N} \times \sum_{i=1}^{N} x_i$$
  
  *na.rm* (boolean : true or false) = remove NA values ?

- Standard-deviation $\sigma$ : *sd* function
  Parameters :    x = list of continuous values
  
  $$\sigma = \sqrt{\frac{1}{N} \times \sum_{i=1}^{N} (x_i - \hat{x})^2}$$
  
  *na.rm* (boolean : true or false) = remove NA values ?

- Median : *median* function
  Parameters :    x = list of continuous values
                  *na.rm* (boolean : true or false) = remove NA values ?

# DESCRIPTIVE STATISTICS
## FUNCTIONS FOR CONTINUOUS VARIABLES

- Mode (available in package *DescTools*) : *mode* function (most frequent value)
  Parameters :     *x* = list of continuous or categorical values
                   *na.rm* (boolean : true or false) = remove NA values ?

- Minimum / Maximum : *min* and *max* functions (range = max – min)
  Parameters :     *x* = list of continuous values
                   *na.rm* (boolean : true or false) = remove NA values ?

- Quartiles (Q1, Q2 (median), Q3 and more) : *quantiles* function
  (interquartile range also called IQR = Q3 – Q1)
  Parameters :     *x* = list of continuous values
                   *probs* = list of probabilities (Q1 : 0.25, Q2 : 0.5 and Q3 : 0.75)
                   *na.rm* (Boolean : true or false) = remove NA values ?

# DESCRIPTIVE STATISTICS
## FUNCTIONS FOR CONTINUOUS VARIABLES

- **Correlation between two variables** (available in package *stats*) : *cor* function
  Parameters :     *x* = list of continuous values
             *y* = list of continuous values
             *method* (*pearson*, *spearman*) = in case of non-normality of data : use spearman method.

Values : from **-1** (perfect **negative** correlation) to **+1** (perfect **positive** correlation)

- **Covariance between two variables** (available in package *stats*) : *cov* function
  Parameters :     *x* = list of continuous values
             *y* = list of continuous values
             *method* (*pearson*, *spearman*) = in case of non-normality of data : use spearman method.

# DESCRIPTIVE STATISTICS
## FUNCTIONS FOR CATEGORICAL VARIABLES

- Mode (available in package *DescTools*) : *mode* function (most frequent value)
  Parameters :     *x* = list of continuous or categorical values
                   *na.rm* (boolean : true or false) = remove NA values ?


- Contingency table : use *table* function :
  Parameters :     *x* = dataset (two categorical variables)
                   *useNA* ("*no*", "*ifany*" or "*always*") = remove NA values ?


- Percentages :     use a **contingency table** to calculate % with *sum* function
                   **percentages within groups** : use *group_by* function

# DESCRIPTIVE STATISTICS
## DESCRIPTIVE TABLE IN SCIENTIFIC PAPERS

- Descriptive table is **mandatory** in research papers

- Goal : give to reviewers **an overview** of the data used

- Content of table :
  - ➤ **Normally** distributed parameter : mean ± standard-deviation
  - ➤ **Not-normally** distributed parameter : median (quartile 1 ; quartile 3)
  - ➤ **Categorical** parameter : count (%) for each modality

Table 1. Demographics

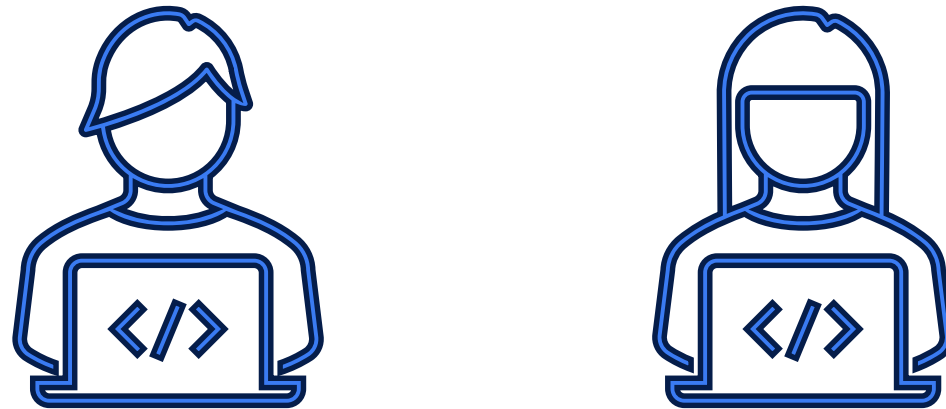|  | Not diabetic (N=9,850) | Diabetic (N=499) |
|---|---|---|
| Age (years) | 46.92 (17.19) | 60.69 (11.47) |
| Weight (kg) | 71.66 (15.22) | 76.67 (17.18) |
| Systolic blood pressure | 130.09 (22.76) | 146.65 (28.39) |
| Sex |  |  |
| Male | 4,698 (47.7%) | 217 (43.5%) |
| Female | 5,152 (52.3%) | 282 (56.5%) |
| Race |  |  |
| White | 8,659 (87.9%) | 404 (81.0%) |
| Black | 1,000 (10.2%) | 86 (17.2%) |
| Other | 191 (1.9%) | 9 (1.8%) |

Total sample: N = 10,349

# DESCRIPTIVE STATISTICS



Live demo

# DESCRIPTIVE STATISTICS

Time to play !
(30 minutes)

# QUESTIONS

07

THANK
YOU
FOR
YOUR
ATTENTION

SEPTEMBER 2025