# One step back to get inertia

## Hardware

## Software



| JupyterLab | Google Colab | Zeppelin |
|---|---|---|
| JetBrains Datalore | Kaggle | Mode Notebooks |
| Observable | Databricks notebooks | Visual Studio Code |
| Amazon SageMaker | CoCalc | Hex |
| Nextjournal | DataCamp Workspace | Deepnote |

## Data

| Kaggle | Google Dataset Search | DrivenData |
|---|---|---|
| UC Irvine Machine Learni... | Dataworld | Data.gov |
| Topcoder | GitHub | FiveThirtyEight |
| InnoCentive | Awesome-public datasets | Zindi |
| Data hub | KDnuggets | Tianchi |
| World Bank | Open data portal | Dataportals |
| HackerEarth | HackerRank | CodaLab |

# Exploratory Data Analysis

The Essence of Exploratory Data Analysis

# GARBAGE IN, GARBAGE OUT.

It is crucial to invest time and effort in **data cleaning, preprocessing, and validation** to minimize the impact of "garbage" data on the overall analysis or decision-making process.
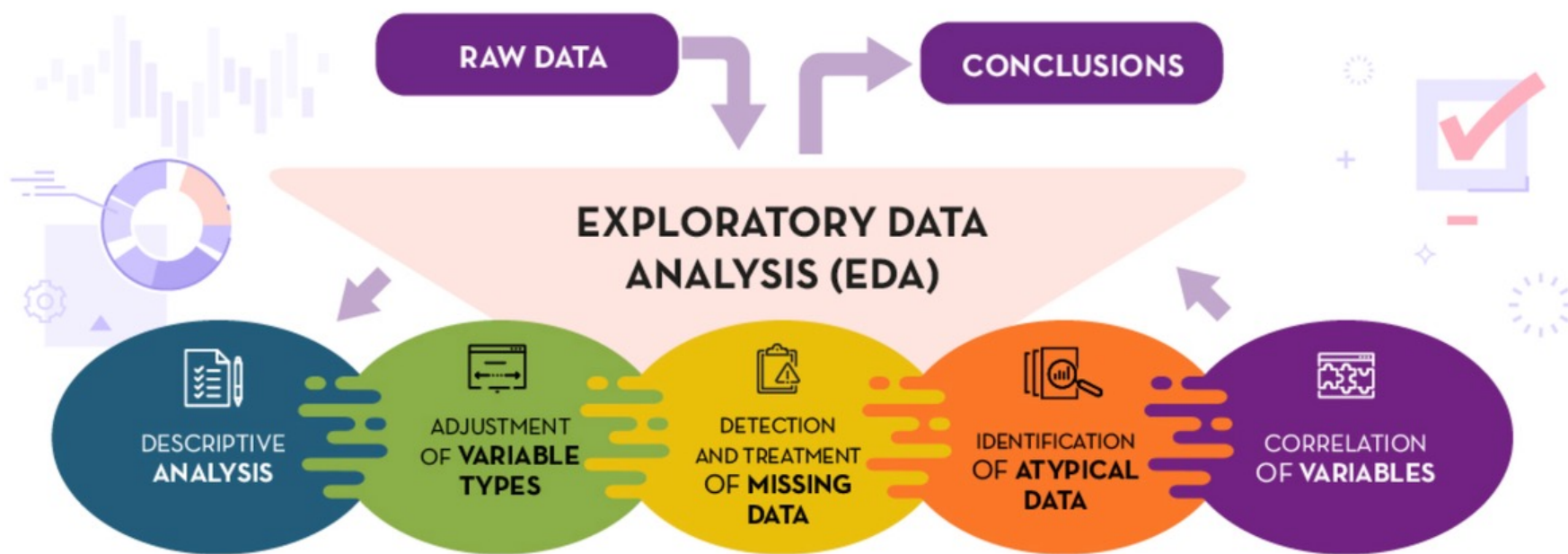
## Exploratory Data Analysis



**RAW DATA** → **CONCLUSIONS**

**EXPLORATORY DATA ANALYSIS (EDA)**

- DESCRIPTIVE ANALYSIS
- ADJUSTMENT OF VARIABLE TYPES
- DETECTION AND TREATMENT OF MISSING DATA
- IDENTIFICATION OF ATYPICAL DATA
- CORRELATION OF VARIABLES

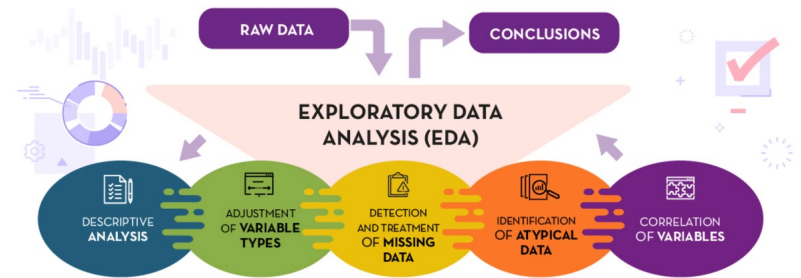Image Credit: https://datos.gob.es/en/documentacion/practical-introductory-guide-exploratory-data-analysis

## o. Methodology



- A practical case

- R / Python / Excel / …
- Not efficient code, but illustrative
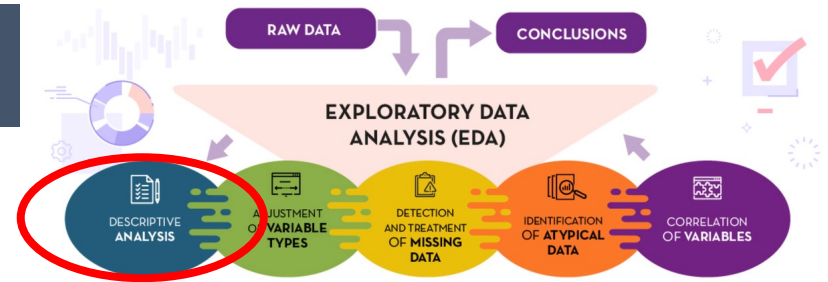- Database:

    **url** =
'https://github.com/DIAGNijmegen/picai_labels/blob/main/clinical_information/marksheet.csv?raw=true'

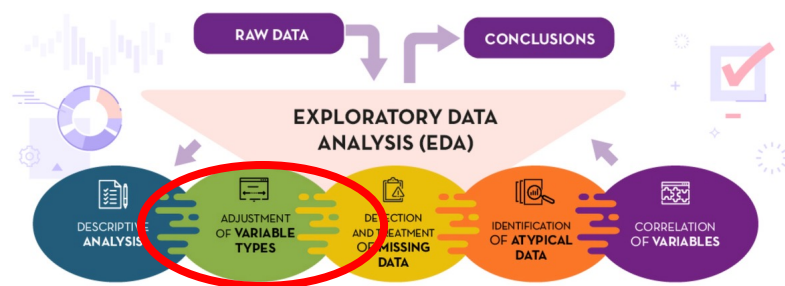- Import libraries and packages

# Exploratory Data Analysis

## 1. Descriptive Analysis

- We will apply **descriptive statistics functions** to explore the structure of the data set and examine the data and variables it presents.

- It will be very useful to use **certain graphic representations** that will help you understand the shape of the data distributions.
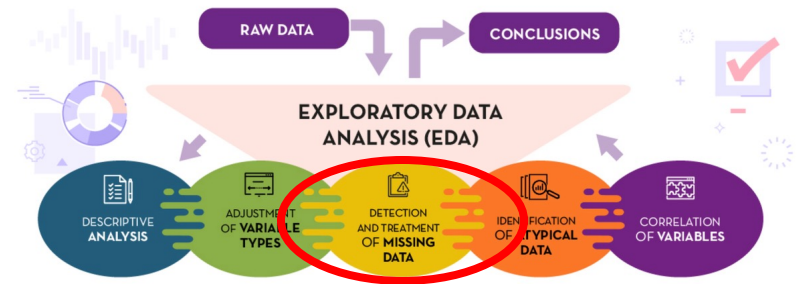
# 2. Variable Types

- After loading the data into the work environment, you must verify that each variable has been stored with the **corresponding value type**.

- Usual types of tabular variables:
  - **numeric**:  stores numbers that can be float or integers.
  - **character**: holds text strings.
  - **categorical**: contains a limited number of categories.
  - **logical** or boolean: binary variables (TRUE/FALSE or 0/1).
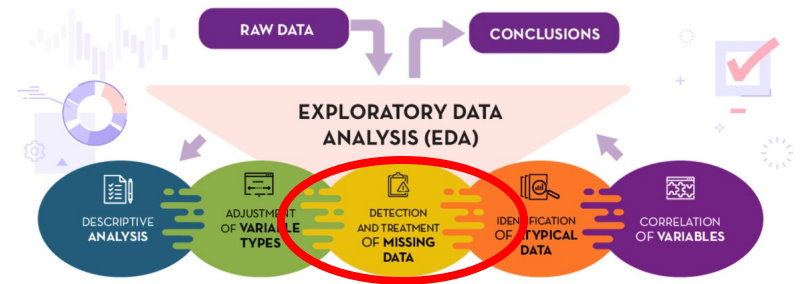  - **date**: stores specific time intervals.

## 3. Missing Data



- Dealing with data sets in which there are **missing data** can cause **problems when applying different statistical analysis or generating graphical representations**.

- In order to avoid future problems, it is necessary to learn how to detect and apply some type of **treatment**.
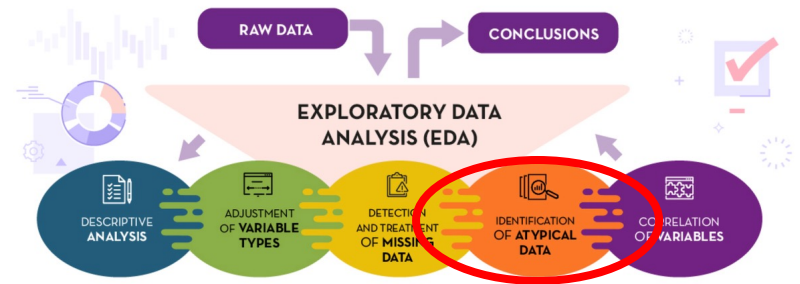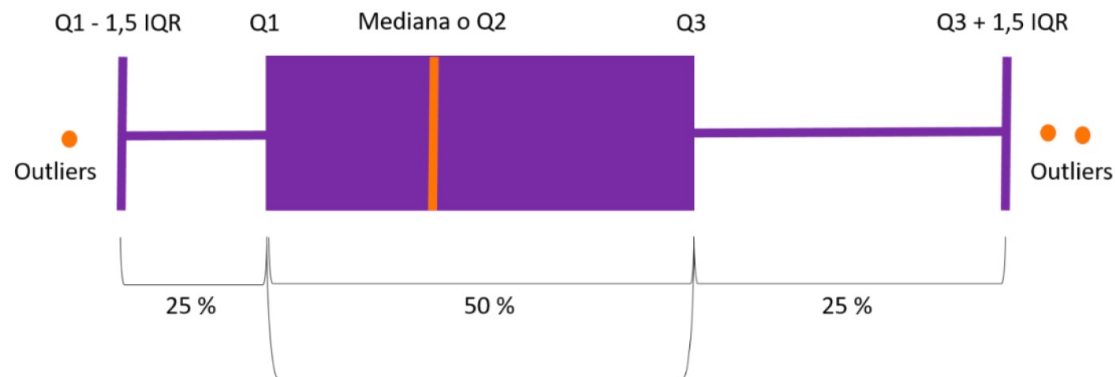
# 3. Missing Data

- There are several ways to deal with missing values:
  - **Fill in** the values with the mean, median or the most frequent value of the variable.
  - **Fill in** missing values with the value directly before or after them in the row or column.
  - **Replace** all missing data with **o**, if they are numerical values.
  - **Delete rows** that have missing values, as long as the data set is large enough.
  - **Eliminate variables** that present a percentage greater than 50% of missing data.
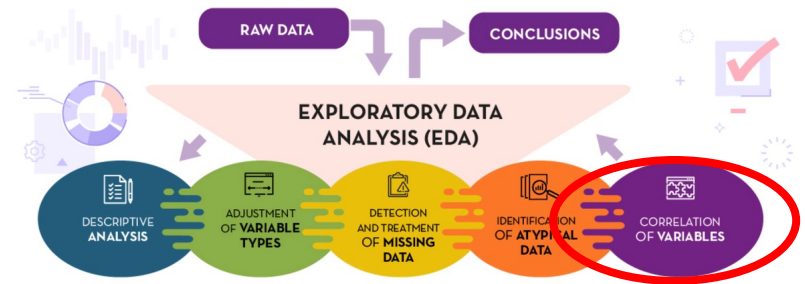
9

## 4. Atypical Data. Outliers



- An **outlier** value is an observation that is significantly different from the rest of the data presented by a variable, of such magnitude that it can be considered an **anomalous value**.

- **IQR-based** or **Percentile** methods.

## 5. Correlation Analysis



- **Correlation** determines the linear relationship between two or more variables, that is, the strength and direction of a possible relationship between variables.