OPEN DATA ANALYSIS WITH OPEN SOURCE

TOOLS

(Part 1)



WHAT IS IT?

Data analysis is a systematic process that allows transforming raw information into valuable knowledge, providing:

- ► Informed decision-making.
- ► Identification of patterns and trends.
- ► In-depth understanding of complex data sets.





///////

Data cleansing Data conversion Data analysis

1. Data preparation and debugging



Objetive: to transform raw data into clean, structured sets.



Examples of tools



OpenRefine Data cleansing and transformation. It offers an intuitive graphical interface that is cross-platform (requires Java). Detectability:

- Duplicated
- · Incomplete data
- Structural inconsistencies



Talend Open Studio ETL (Extract, Transform, Load) tool. (!) Requires intermediate

programming skills. Permitted:

- Component programming
- Integrating data from multiple sources

2. Data conversion



Objetive: to adapt the format of the data to facilitate its analysis.



Examples of tools



COMPLETELY Mr Data Converter

- · Conversion between CSV, Excel, JSON, HTM, XLM formats
- Simple web interface
- · No installation required



COMPLETELY Pandoc

- Universal document conversion
- Supports more than 20 different formats
- · Powerful command line



COMPLETELY **Tabula**

- · Extraction of tables from PDF
- · Convert documents into reusable formats
- · Useful for reports and official documentation



3. Data analysis



Objetive: to explore, process and gain insights from datasets.



User-friendly analysis software



COMPLETELY WEKA

- · Machine learning and data mining
- Graphical interface · Integration with scikitlearn, R and Deeplearning
- · Ideal for beginners in machine learning



COMPLETELY **ORANGE**

- Paradigm drag and drop (drag and drop)
- Interactive visualisations
- · Accessible statistical analysis



FREE IN BASIC **KNIME**

- · Visual data analysis
- Workflows by connecting nodes • Extensive library of components



Development environments



Jupyter Notebook



COMPLETELY

- Combination of code, visualisations and narrative
- Supports multiple languages · Ideal for reproducibility



FREE IN BASIC RStudio

- Complete R language environment · Console, editor and visualisation
- integration Advanced statistical tools



Programming languages for data analysis COMPLETELY



Especialist in statistics • Powerful for statistical analysis

- and visualisation
 - Featured libraries: » Tidyverse
 - » ggplot2



COMPLETELY **Python**

A versatile language

- Recommendation: use Anaconda for environment management
- The main ones <u>libraries for analysis</u> are: » Pandas (data manipulation)
 - » NumPy (numerical calculation)
 - » scikit-learn (machine learning) » Matplotlib (visualisation)



Streamlit

Emerging tools



· Rapid creation of data-driven web

Apache Spark

• For Big data.

applications Requires Python only

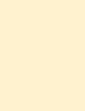
COMPLETELY

 Rapid prototyping of dashboards. » In this <u>practical exercise</u> we use

· Distributed processing

• APIs in Python, R and Scala

it to create a public data chat COMPLETELY



High performance

Polars

• Optimised pandas alternative Parallel processing

COMPLETELY





Recommendation Start with simple tools such as Jupyter and Python, and gradually explore more advanced options as your analysis needs change.

Find out here the benefits and steps of







