# Data Analysis with Python

## Full tutorial for beginners

IMPORTANT: We're in the process of adapting the projects into DataWars.io for your simplicity.

(stay tuned 👍)

# About this tutorial

# What is Data Analysis?

# What is Data Analysis

> *A process of inspecting, cleansing, transforming and modeling data with the goal of discovering useful information, informing conclusion and supporting decision-making.*

[Definition by Wikipedia](#).

RMOTR
BY INE

# What is Data Analysis

> *A process of* ***inspecting, cleansing, transforming*** *and modeling data with the goal of discovering useful information, informing conclusion and supporting decision-making.*

Definition by Wikipedia.

# What is Data Analysis

*> A process of inspecting, cleansing, transforming and **modeling data** with the goal of discovering useful information, informing conclusion and supporting decision-making.*

Definition by Wikipedia.

RMOTR BY INE

# What is Data Analysis

*> A process of inspecting, cleansing, transforming and modeling data with the goal of **discovering useful information**, informing conclusion and supporting decision-making.*

Definition by Wikipedia.

RMOTR
BY INE

# What is Data Analysis

> *A process of inspecting, cleansing, transforming and modeling data with the goal of discovering useful information,* ***informing conclusion and supporting decision-making****.*

Definition by Wikipedia.

RMOTR BY INE

Data Analysis Tools

# Auto-managed closed tools

# Programming Languages

## Auto-managed closed tools

👎 Closed Source 🙅‍♂️

👎 Expensive 💸

👎 Limited 😩

👍 Easy to learn 👩‍💻

## Programming Languages

👍 Open Source 🤩

👍 Free (or very cheap) 🤑

👎 Extremely Powerful 💪

👎 Steep learning curve 👩‍💻

# Why Python for Data Analysis?

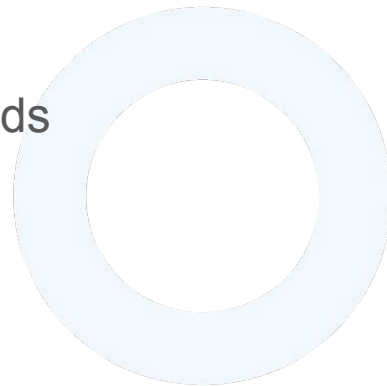# Why Python for Data Analysis?

*Why would we choose Python over R or Julia?*

👍 very simple and intuitive to learn

👍 "correct" language

👍 powerful libraries (not just for Data Analysis)

👍 free and open source

👍 amazing community, docs and conferences

**RMOTR** BY INE

# When to choose R?

*Python, sadly, is not always the answer*

- When R Studio is needed

- When dealing with advanced statistical methods

- When extreme performance is needed

**RMOTR** BY INE

# The Data Analysis Process

| Data Extraction | Data Cleaning | Data Wrangling | Analysis | Action |
|---|---|---|---|---|

**Data Extraction**
- SQL
- Scrapping
- File Formats
  - CSV
  - JSON
  - XML
- Consulting APIs
- Buying Data
- Distributed Databases

**Data Cleaning**
- Missing values and empty data
- Data imputation
- Incorrect types
- Incorrect or invalid values
- Outliers and non relevant data
- Statistical sanitization

**Data Wrangling**
- Hierarchical Data
- Handling categorical data
- Reshaping and transforming structures
- Indexing data for quick access
- Merging, combining and joining data

**Analysis**
- Exploration
- Building statistical models
- Visualization and representations
- Correlation vs Causation analysis
- Hypothesis testing
- Statistical analysis
- Reporting

**Action**
- Building Machine Learning Models
- Feature Engineering
- Moving ML into production
- Building ETL pipelines
- Live dashboard and reporting
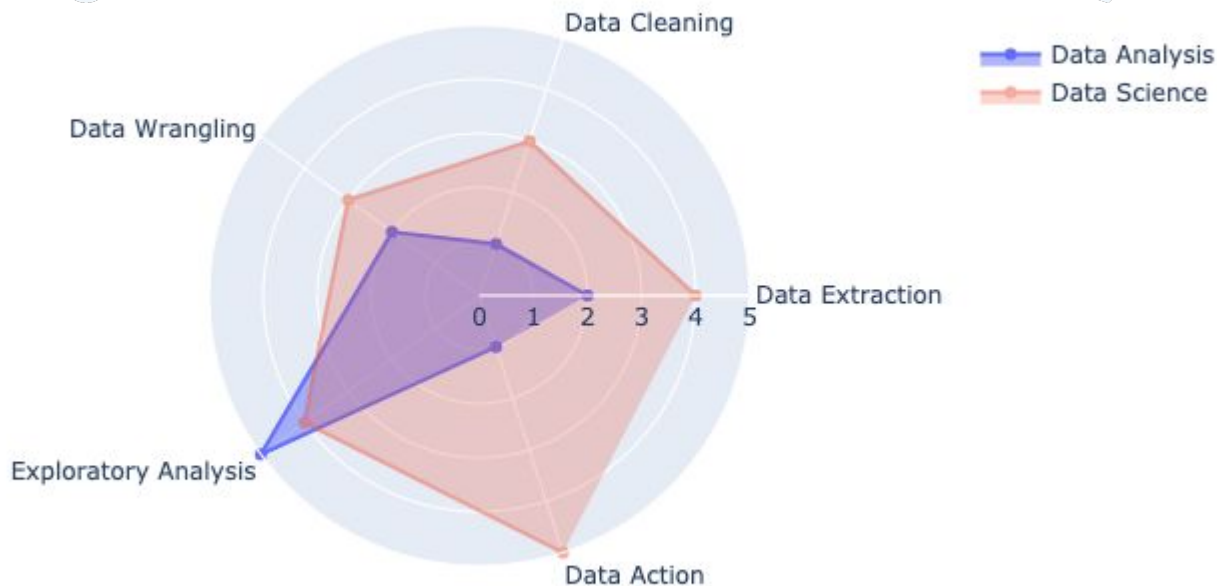- Decision making and real-life tests

# Data Analysis
# Vs
# Data Science

DATA ANALYSIS VS DATA SCIENCE
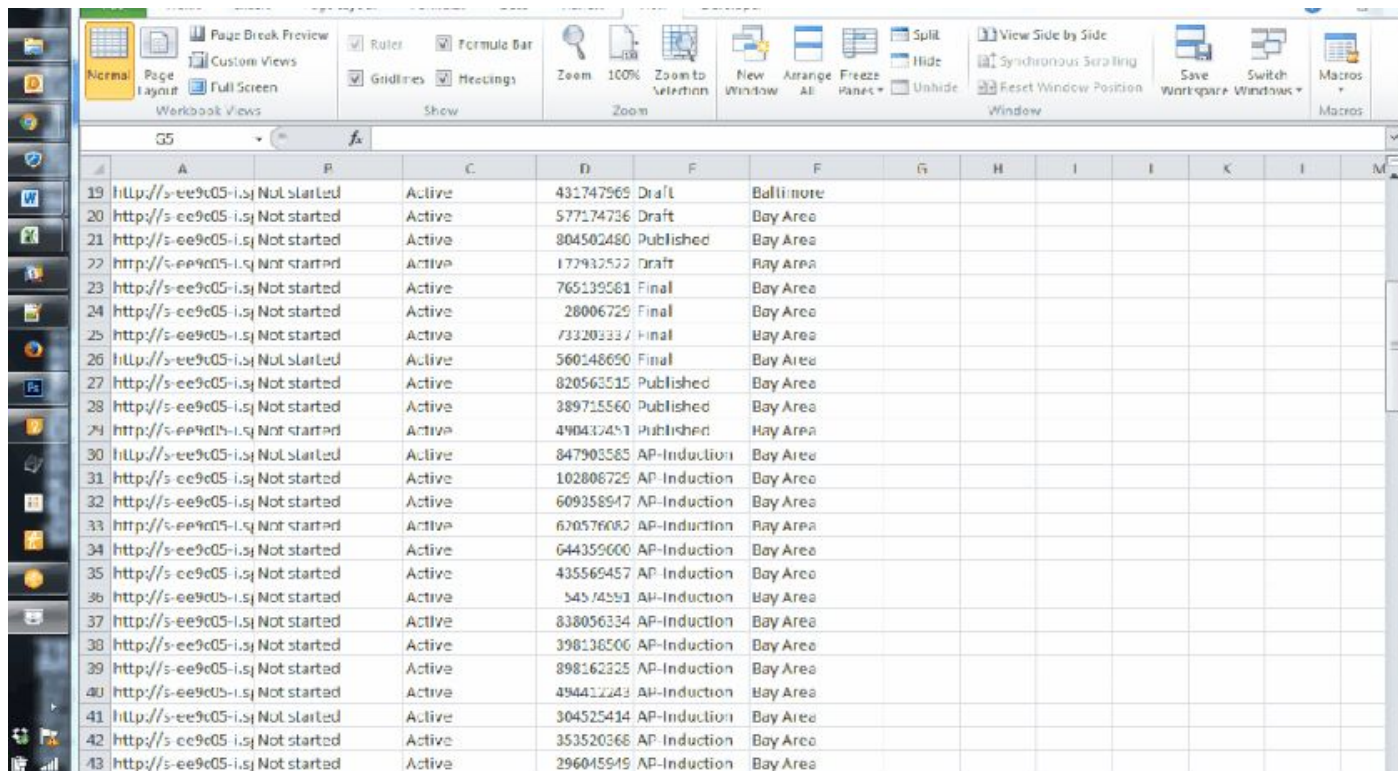
# The traditional view

# Python & PyData Ecosystem

# The libraries we use...

- **pandas**: The cornerstone of our Data Analysis job with Python

- **matplotlib**: The foundational library for visualizations. Other libraries we'll use will be built on top of matplotlib.

- **numpy**: The numeric library that serves as the foundation of all calculations in Python.

- **seaborn**: A statistical visualization tool built on top of matplotlib.

- **statsmodels**: A library with many advanced statistical functions.

- **scipy**: Advanced scientific computing, including functions for optimization, linear algebra, image processing and much more.

- **scikit-learn**: The most popular machine learning library for Python (not deep learning)

# They're all visual tools...

# Thinking like a
# Python Data Analyst

# And finally, why Python?

# >20%

Salary increase for a Data Analyst that knows Python and SQL.

# About this tutorial

1.  What is Data Analysis

2.  **Real Example Data Analysis with Python**

3.  How to use Jupyter Notebooks

4.  Intro to NumPy (exercises included)

5.  Intro to Pandas (exercises included)

6.  Data Cleaning

7.  Reading Data SQL, CSVs, APIs, etc

8.  Python in Under 10 Minutes