



# DATA ANALYSIS WITH PYTHON

# TUTORIAL



Introduction to Data Analysis



NumPy

# I. INTRODUCTION TO DATA ANALYSIS – TUTORIAL

1. What is Data Analysis
2. Example Data Analysis with Python
3. How to use Jupyter Notebooks
4. Intro to Numpy
5. Intro to Pandas
6. Data Cleaning
7. Reading Data SQL, CSVs, APIs, etc.
8. Python

## ■ 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”

# WHY PYTHON FOR DATA ANALYSIS?

- Very simple and intuitive to learn
- “Correct” language
- Powerful libraries (not just for Data Analysis)
- Free and open source
- Amazing community, docs and conferences



## When to choose R?

- When R Studio is needed
- When dealing with advanced statistical methods
- When extreme performance is needed



# THE DATA ANALYSIS PROCESS

## 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

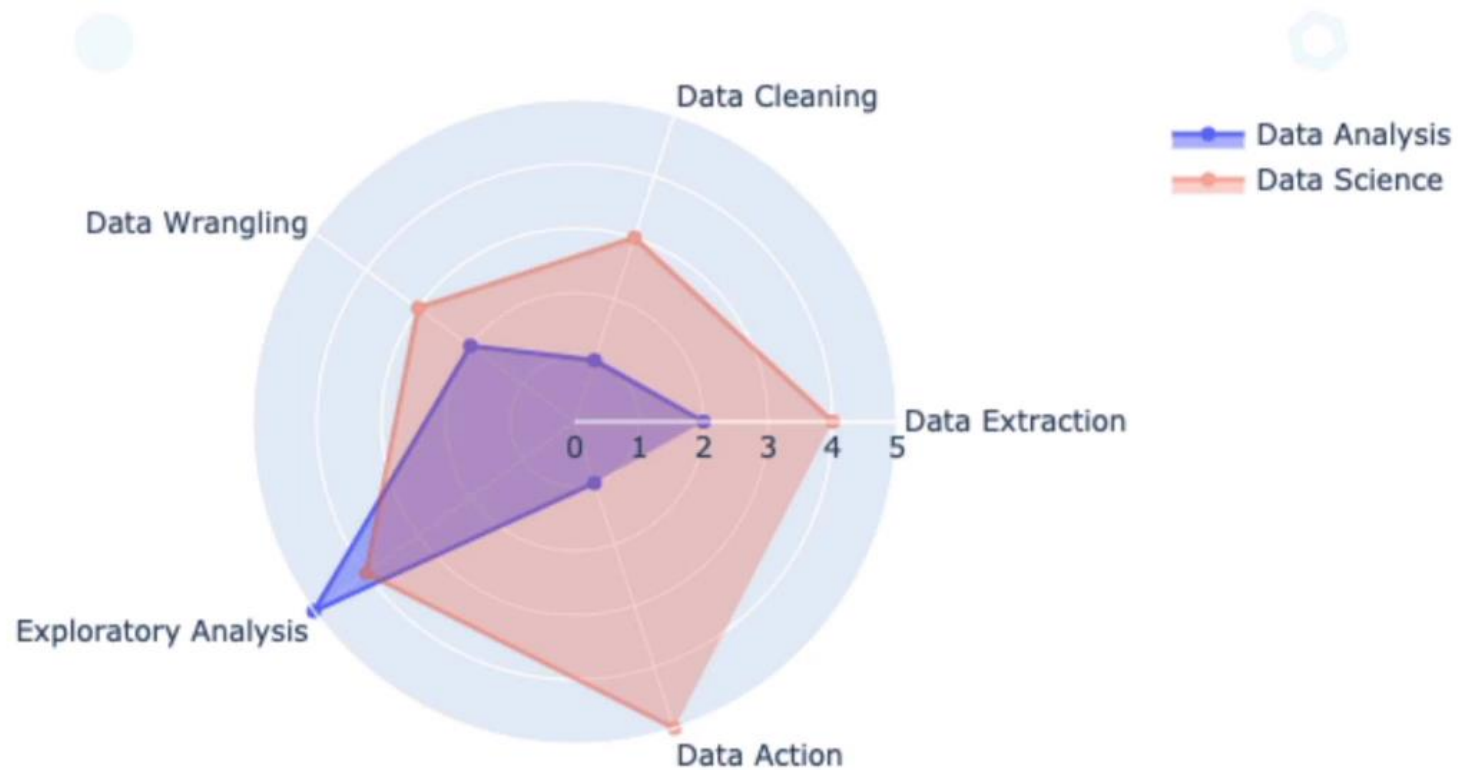
- Exploration
- Building statistical models
- Visualization and representations
- Correlation vs Causation analysis
- Hypothesis testing
- Statistical analysis
- Reporting

## 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

## The traditional view

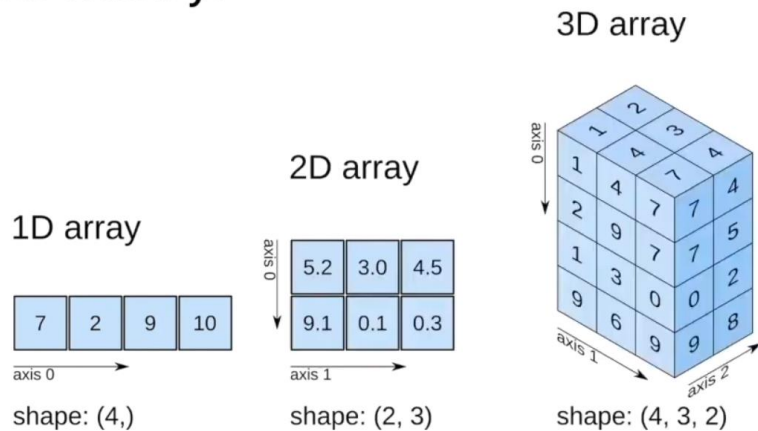


# PYTHON & PYTHON DATA ECOSYSTEM

- [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)

## 2. NUMPY

### What is NumPy?

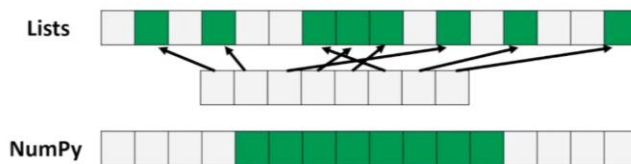


### How are Lists different from Numpy?



- Faster to read less bytes of memory.
- No type checking when iterating through objects.

### Why is NumPy Faster? - Contiguous Memory



Benefits:

- SIMD Vector Processing
- Effective Cache Utilization

### Why is NumPy Faster? - Fixed Type

5 → Binary 00000101

Int32  
00000000 00000000 00000000 00000101

NumPy

Lists

00000000 00000000 00000000 00011100  
00000001 00111101 11111110 10111100 00011010 11011101 10100100 11011000  
11001010 10111110 01100001 01000100 11111100 00000000 11001100 01011111  
00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000101

Size  
Reference Count  
Object Type  
Object Value

Lists



# NUMPY

## How are Lists different from Numpy?

### Lists

```
a = [1,3,5]
```

```
b = [1,2,3]
```

```
a*b = ERROR
```

### NumPy

```
a = np.array([1,3,5])
```

```
b = np.array([1,2,3])
```

```
a*b = np.array([1,6,15])
```

## Applications of NumPy?

Mathematics (MATLAB Replacement)

Plotting (Matplotlib)

Backend (Pandas, Connect 4, Digital Photography)

Machine Learning



THANKS