# Data Management & Visualization

Eng- Mohamed Khaled Idris
Eng- Mayar Swilam

# Review:



Step -1
Collection of Data from Various source

Step -2
Data cleaning and Feature Engineering

Step -3
Model building for selecting correct ML Algorithm

Step -4
Evaluate Model

Step -5
Model Deployment

2. **Visualization** Libraries

**Matplotlib**
(plots & graphs, most popular)

**Seaborn**
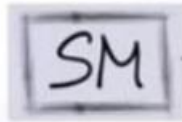(plots : heat maps, time series, violin plots)

3. **Algorithmic** libraries

**Scikit-learn**
(Machine Learning : regression, classification,... )

**Statsmodels**
(Explore data, estimate statistical models, and perform statistical tests.)

**TENSORFLOW & PYTORCH**

# Data Management

- In order to manage data through python, there are many ways to do it !
- Also, you may need to do some mathematical operations to matrices, lists or any structure
- You can simply manage your data as the same as SQL

## HOW THEN ?

# What is NumPy ?

- NumPy is an open-source Python library that facilitates efficient numerical operations on large quantities of data.

- There are a few functions that exist in NumPy that we use on pandas DataFrames.

- For us, the most important part about NumPy is that pandas is built on top of it. So, NumPy is a dependency of Pandas.
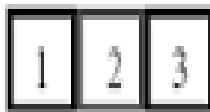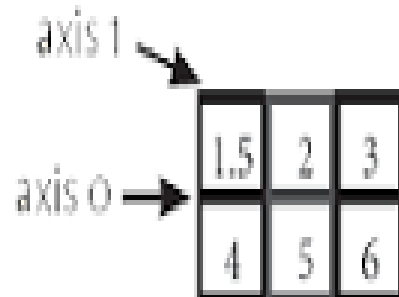
# NumPy Arrays

- **NumPy arrays are unique in that they are more flexible than normal Python lists.**

- **They are called ndarrays since they can have any number (n) of dimensions (d).**

- **They hold a collection of items of any one data type and can be either a vector (one-dimensional) or a matrix (multi-dimensional).**

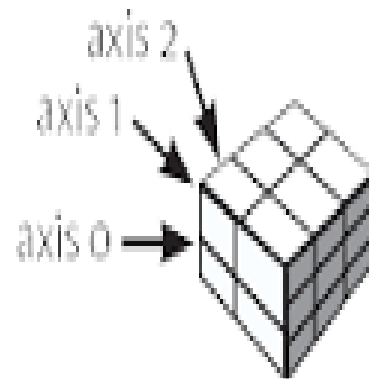- **NumPy arrays allow for fast element access and efficient data manipulation.**

# Datacamp Cheat Sheet (NumPy)

## Python For Data Science
## NumPy Cheat Sheet

Learn NumPy online at www.DataCamp.com

### Numpy

The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays

Use the following import convention:

```
>>> import numpy as np
```

#### NumPy Arrays

1D array    2D array    3D array

### Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1.5,2,3), (4,5,6)], dtype = float)
>>> c = np.array([[(1.5,2,3), (4,5,6)],[(3,2,1), (4,5,6)]], dtype = float)
```

#### Initial Placeholders

### Inspecting Your Array

```
>>> a.shape #Array dimensions
>>> len(a) #Length of array
>>> b.ndim #Number of array dimensions
>>> e.size #Number of array elements
>>> b.dtype #Data type of array elements
>>> b.dtype.name #Name of data type
>>> b.astype(int) #Convert an array to a different type
```

### Data Types

```
>>> np.int64 #Signed 64-bit integer types
>>> np.float32 #Standard double-precision floating point
>>> np.complex #Complex numbers represented by 128 floats
>>> np.bool #Boolean type storing TRUE and FALSE values
>>> np.object #Python object type
>>> np.string_ #Fixed-length string type
>>> np.unicode_ #Fixed-length unicode type
```

### Array Mathematics

### Sorting

```
>>> a.sort() #Sort an array
>>> c.sort(axis=0)
```

### Subsetting

# What is Pandas ?

- *Pandas* is a very popular library for working with data (its goal is to be the most powerful and flexible open-source tool, and in our opinion, it has reached that goal).

- DataFrames are at the center of pandas. A DataFrame is structured like a table or spreadsheet. The rows and the columns both have indexes, and you can perform operations on rows or columns separately.

# DataFrame

- A pandas DataFrame can be easily changed and manipulated. Pandas has helpful functions for handling missing data, performing operations on columns and rows, and transforming data.

- If that wasn't enough, a lot of SQL functions have counterparts in pandas, such as join, merge, filter by, and group by.

- With all of these powerful tools, it should come as no surprise that pandas is very popular among data scientists.

# Pandas

The **Pandas** library is built on NumPy and provides easy-to-use **data structures** and **data analysis** tools for the Python programming language.

Use the following import convention:

```
>>> import pandas as pd
```

## > Pandas Data Structures

### Series

A one-dimensional labeled array capable of holding any data type

| Index → | a | 3 |
| | b | -5 |
| | c | |

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

### Dataframe

A two-dimensional labeled data structure with columns of potentially different types

| Columns → | Country | Capital | Population |
| Index → 0 | Belgium | Brussels | 11190846 |
| 1 | India | New Delhi | 1303171035 |
| 2 | Brazil | Brasília | 207847528 |

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
            'Capital': ['Brussels', 'New Delhi', 'Brasília'],
            'Population': [11190846, 1303171035, 207847528]}
>>> df = pd.DataFrame(data,
              columns=['Country', 'Capital', 'Population'])
```

## > Dropping

```
>>> s.drop(['a', 'c'])   #Drop values from rows (axis=0)
>>> df.drop('Country', axis=1)   #Drop values from columns(axis=1)
```

## > Asking For Help

```
>>> help(pd.Series.loc)
```

## > Sort & Rank

```
>>> df.sort_index()   #Sort by labels along an axis
>>> df.sort_values(by='Country')   #Sort by the values along an axis
>>> df.rank()   #Assign ranks to entries
```

## Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')
>>> df.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')
```

**Read multiple sheets from the same file**

```
>>> xlsx = pd.ExcelFile('file.xls')
>>> df = pd.read_excel(xlsx, 'Sheet1')
```

## Read and Write to SQL Query or Database Table

```
>>> from sqlalchemy import create_engine
>>> engine = create_engine('sqlite:///:memory:')
>>> pd.read_sql("SELECT * FROM my_table;", engine)
>>> pd.read_sql_table('my_table', engine)
>>> pd.read_sql_query("SELECT * FROM my_table;", engine)
```

read_sql() is a convenience wrapper around read_sql_table() and read_sql_query()

```
>>> df.to_sql('myDF', engine)
```

## > Selection                    Also see NumPy Arrays

### Getting

```
>>> df.iloc[2]   #Select single row of subset of rows
Country       Brazil
Capital       Brasília
Population    207847528
>>> df.loc[:,'Capital']   #Select a single column of subset of columns
0    Brussels
1    New Delhi
2    Brasília
>>> df.loc[1,'Capital']   #Select rows and columns
'New Delhi'
```

**Boolean Indexing**

```
>>> s[~(s > 1)]   #Series s where value is not >1
>>> s[(s < -1) | (s > 2)]   #s where value is <-1 or >2
>>> df[df['Population']>1200000000]   #Use filter to adjust DataFrame
```

**Setting**

```
>>> s['a'] = 6   #Set index a of Series s to 6
```

## Summary

```
>>> df.sum()   #Sum of values
>>> df.cumsum()   #Cumulative sum of values
>>> df.min()/df.max()   #Minimum/maximum values
>>> df.idxmin()/df.idxmax()   #Minimum/Maximum index value
>>> df.describe()   #Summary statistics
>>> df.mean()   #Mean of values
>>> df.median()   #Median of values
```

## > Applying Functions

```
>>> f = lambda x: x*2
>>> df.apply(f)   #Apply function
>>> df.applymap(f)   #Apply function element-wise
```

## > Data Alignment

### Internal Data Alignment

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
```

### Arithmetic Operations with Fill Methods

You can also do the alignment yourself with the help of the fill methods:

- A picture is worth a thousand words

- For exploratory data analysis

- Communicate data clearly

- Share unbiased representation of data

- Support recommendations to different stakeholders

# Why Visualization ?
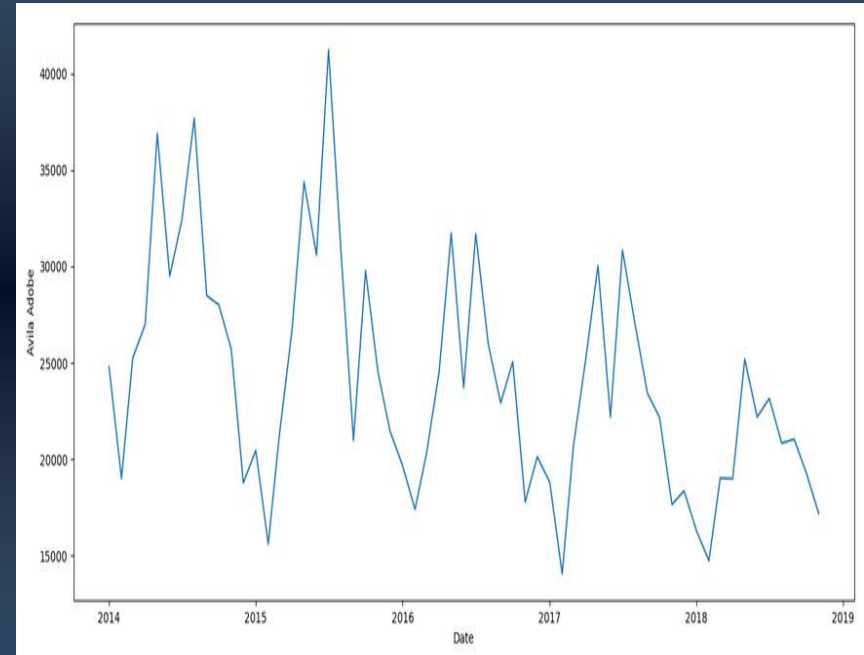
# The important of Visualization

- Considering only the years represented in the figure, which countries spent at least 5 consecutive years in the #1 ranked spot?
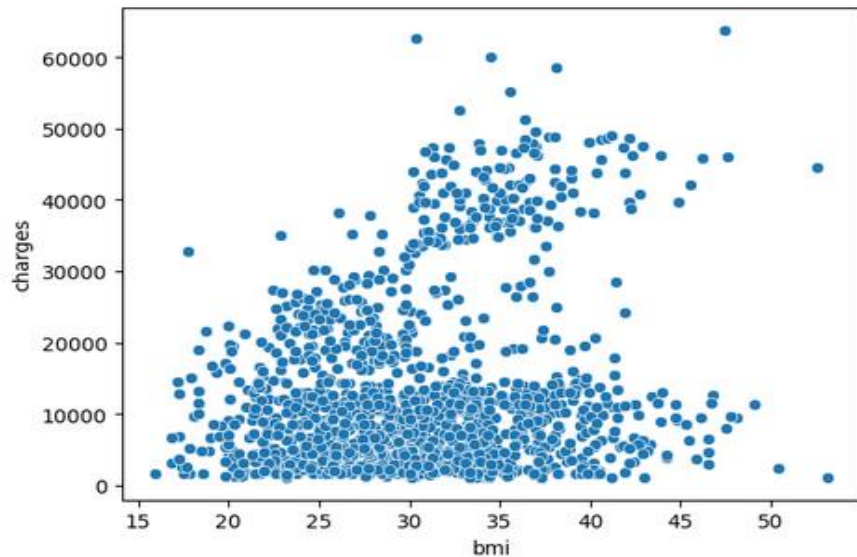
Does Avila Adobe get more visitors:

- **in September-February (in LA, the fall and winter months), or**

- **in March-August (in LA, the spring and summer)?**

Using this information, when should the museum staff additional seasonal employees?

IS THERE ANY RELATION BETWEEN BODY MASS INDEX AND INSURANCE CHARGES ?

# Matplotlib

- Neurobiologist
- Part of a team analyzing **Electrocorticography Signals (ECoG)**
  - **Electrocorticography** is the process of recording electrical activity in the brain

- The team
  - used a proprietary software (**MATLAB** based version) for analysis
  - had only one license and were taking turns in using it

- **John** replace the proprietary software with **Matplotlib**

## Matplotlib

Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.

### > Prepare The Data

#### 1D Data

#### 2D Data or Images

### > Create Plot

#### Figure

#### Axes

#### 2D Data

#### Vector Fields

#### Data Distributions

#### Plot Anatomy

#### Workflow

#### Markers

#### Linestyles

#### Text & Annotations

#### Mathtext

#### Limits, Legends

# Datacamp Cheat Sheet (Matplotlib)

https://images.datacamp.com/image/upload/v1676360378/Marketing/Blog/Matplotlib_Cheat_Sheet.pdf

Questions

# Thank You