

## Data Science ?

- Doing science like studies on data
- Involves running experiments, doing analysis, coming to conclusions
- Involves building models and testing them.
- like cleaning a beaker before testing for chemicals, we clean data.

## How to do Data Science ?

- Involves a lot of steps like data fetching, process cleaning, analyzing, presenting, ML modelling, deployment, scaling etc.

solve problems like :

- ETL ( Extract Load Transform) jobs
- A/B Testing
- Data scraping and cleaning
- Doing analyses and presenting results
- Building and maintaining dashboards
- Building and deploying ML models

## Numpy: Matrix & Tensor Operations

## Pandas: Reading, cleaning & operating on Data

## opencv: Reading & basic operations on Images

## Seaborn & Matplotlib: Plotting of data

## Scrapy & beautifulSoup: Fetching data from the internet

## My.SQL: Fetching data from Database

Numpy -

what is numpy?

- Fundamental of every data science library.
- Almost every data science library in python is itself written in Numpy.
- Numpy deals with fundamental building block of all data science algos: vectors and Matrices.
- If data science was a car, Numpy would be the wheels.

15 when to use Numpy?

- Optimized for doing numerical data science work.
- Numpy is a python library written in C, making it blazingly fast.
- It has a lot of submodules capable of doing a variety of tasks related to data science.

When NOT to use Numpy?

- Doesn't work well with Non-numerical data.
- If a lot of data is in text format, like review, you should not directly use Numpy.
- Numpy always runs on your CPU. If however you are doing deep learning work on a GPU, Numpy won't work.

30 what is a vector?

- List of numbers.

e.g.: [1, 2, 3, 4, 10, -5]  
 ↑  
 int

[1, 'A', 2, 'B', 3] → LIST  
 ↑  
 T... int

what is a matrix?

List of vectors i.e. a list of lists of numbers  
e.g -  $\begin{bmatrix} [1, 2, 3] \\ [4, 5, 6] \end{bmatrix}$

### Jupyter

↳ pip install jupyter

10 jupyter notebook

new → python

### 15 Some practice

In [1]: import numpy as np

In [2]: python\_list = [1, 2, 3, 4, -5]

In [3]: python\_list

Out [3]: [1, 2, 3, 4, -5]

In [4]: np.array → to convert list to vector

Out [4]: array ([1, 2, 3, 4, -5]) → 1D vector

In [5]: numpy\_arr = np.array (python\_list)

In [6]: numpy\_arr.shape

Out [6]: (5,) → tuple

In [7]: len (numpy\_arr.shape)

Out [7]: 1

In [8]: python\_list\_list = [[1, 2], [3, 4]]

In [9]: python\_list\_list

Out [9]: [[1, 2], [3, 4]]

In [10]: np.array (python\_list\_list)

Out [10]: array ([[1, 2], [3, 4]]) → 2D list

Numpy : Data Types

Numpy arrays have data types.

- Unlike in Python, numpy arrays have a specific data types
- All the elements of an array must be same data type.
- Numpy also has a few unique things to handle non-numeric data, infinities etc.

15 Data types :

int	(Integers)	(int 8, int 16, int 32, int 64)	$2^7 = 128$
uint	(Unsigned Integers)	(uint 8, uint 16, uint 32, uint 64)	$2^{16} = 65536$
bool	(Boolean (True or False))		$2^1 = 2$
float	(Floating-point numbers)	(float 16, float 32, float 64)	$2^{-149} \text{ to } 2^{148}$
complex	(Complex-values floating-point numbers)		$2^{-149} \text{ to } 2^{148}$

Numpy : Slicing

- Used for manipulating arrays
- Reduce length selection, similar to python list
- Can use Boolean arrays for selecting and filtering.

30 arr [:4] (Slice till 4)

(include till index 3)  
as  $3 < 4$ arr [:, :4]  
"select all"

Numpy : Transpose

Swap rows with columns and columns with rows

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}_{3 \times 3} \xrightarrow{\text{Transpose}} \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix}_{3 \times 3}$$

Can change shape of non-square matrix

$$\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{bmatrix}_{3 \times 4} = \begin{bmatrix} 1 & 5 & 9 \\ 2 & 6 & 10 \\ 3 & 7 & 11 \\ 4 & 8 & 12 \end{bmatrix}_{4 \times 3}^T$$

Numpy : Broadcasting

Operations on arrays of different dimensions

Using broadcasting, we can add this vector, and the vector gets broadcasted / copied so that it can be added.

$$\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{bmatrix}_{3 \times 4} + \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{bmatrix}_{3 \times 4} + \begin{bmatrix} 1 & 1 & 1 \\ 2 & 2 & 2 \\ 3 & 3 & 3 \end{bmatrix}_{3 \times 3}$$

same rows ✓

→ Broadcasting

Numpy : Tensors

Like matrix is a 2D array, tensor is a 3D array.

Tensor could be also 4D, 5D ... arrays but mostly tensors are 3D arrays.

## Numpy : Tensor Transpose

Transpose is usually done by swapping any 2 axes, or rotating the 3d matrix.

## Pandas

- Pandas is how we read and process data in python.
- Almost every project uses pandas to read data, process data and make it ready for use in ML.
- Uses numpy under the hood to make things fast.
- Similar to using SQL, but more powerful for some use cases.

When do we use Pandas?

- Reading and processing tabular data.
- Reading data from files (like csv) locally or from the cloud. Can also be used to read from SQL.
- Joining multiple data files, removing invalid or wrong data, calculating business relevant numbers and many more things are usually done using pandas.

→ Pandas is also usually used for feature engg for ML use cases.

When NOT use Pandas?

- Non-tabular data should not be read using pandas.
- Images, Not suitable for very large data.

## Numpy : Tensor Transpose

Transpose is usually done by swapping any 2 axes, or rotating the 3d matrix.

## \* Pandas

- Pandas is how we read and process data in python.
- Almost every project uses pandas to read data, process data and make it ready for use in ML.
- Uses numpy under the hood to make things fast.
- Similar to using SQL, but more powerful for some use cases.

### When do we use Pandas?

- Reading and processing tabular data.
- Reading data from files (like csv) locally or from the cloud. Can also be used to read from SQL.
- Joining multiple data files, removing invalid or wrong data, calculating business relevant numbers and many more things are usually done using pandas.
- Pandas is also usually used for feature engg for ML use cases.

### When NOT use Pandas?

- Non - tabular data should not be read using pandas
- Images, Not suitable for very large data.

Dataframe & Series

5 Pandas Tables are called DataFrames and each column for which the DF is made up of is a Series.

Functions like `head`, `tail`, `info` etc. allows to see inside a pandas DataFrame.

10 Pandas can read data from SQL, CSV, Excel, JSON, Pickle etc.

## Data format

## Users

15 CSV (Comma Separated Values) A Plain tabular only format

- XLSX / .XLS

Familiar Microsoft Excel file

• JSON (JavaScript Object Notation) Dictionary-based plain text data format

20

Pickle (Python Object) entire data frame is stored

(.pkl / .pk) as Binary pickle file

## 25 Pandas: Data types

Pandas supports a few more data types than Numpy.

object (str or mixed)

int64

float64

bool

Read data  $\Rightarrow$

import pandas as pd

countries\_data = pd.read\_csv("d/countries.csv")

countries\_data.tail() <sup>last row</sup>

countries\_data["country"].unique()

### Pandas Indexing

- A dataframe index is analogous to python list index.
- In lists, indices are used because they allow very fast & easy access to the list elements. They are simple integers from 0 onwards.
- while usually dataframe indices will be simple integers from 0/1 that is not always the case.
- we can set any column as the index of the dataframe, at which point it will be analogous to keys in a dictionary.
- Dataframe rows can be accessed using loc and iloc.
  - ↳ useful when applied some filtering

### Pandas: Masking

Similar to Numpy we can use Masking to select from a DataFrame.

countries\_data.query("(country == 'Zimbabwe') & (date == '20/11/2011')")

## Pandas: NULL values

5 Data sometimes have blank values, similar to points in yellow.

Pandas often "fill" these blank spaces with a "NaN" or "None".

10 NaN: Stands for "Not a Number". Pandas will often use `NaN` for blank values in numerical columns.

(for strings)  
15 None: Python's equivalent of `NULL`. It will be used to fill blanks for object type columns.

NAT: Stands for "Not a Time". Pandas often use `nat` for blank values in `datetime`

20 however, not all null values will have `NaNs/None` often an invalid value may be used in place of `NaN/None`

Eg: Using `-1` in an `Age` column to show that age is unknown.  
25

## Pandas: Kinds of join/ merge

30 There are 5 main types of joins/ merge

INNER: The default, where only rows which are present on both dataframes are kept.

LEFT: Where rows which are present on both dataframes are kept, along with leftover rows from the left dataframe.

RIGHT: Where rows which are present on both dataframes are kept, along with leftover rows from the right dataframe.

OUTER / FULL: Where rows which are present on both dataframes are kept, along with leftover rows from both left & right dataframes.

ANTI-JOIN: Where rows which are present on both dataframes are removed, keeping only leftover rows from both dataframes.

These are also anti-left and anti-right join which only keep leftover rows from left data frame and right data frame respectively.

### Data Scraping -

Process of getting useful, semi-structured data from the Internet is called Web Scraping.

Basics of HTML, CSS and JSON needed for web scraping.

Types of web scrapping -

- Request/API scraping
- Static HTML scraping
- Dynamic JS scraping

### Request / API scraping -

We directly load data from Internet via an "API". Considered to be the easiest type of scraping.

5

### Static HTML scraping -

We load a simple "static" website. We extract HTML data available on the website. Also considered to be easy.

10

### Dynamic JS scraping -

We load a complex "dynamic" website. We interact with it like a person.

15

#### 1] API Scraping

API (Application Programming Interface) is a way to get an external program to do something.

For scraping, we send a "request" to the API program to get data for us.

JSON is a way to store data in text.

→ We will be using requests python library for easy API scraping.

#### 2] HTML scraping -

Tools - Scrapy & BeautifulSoup.

→ Scrapy will allow us to create a simple way to scrape and crawl our desired website, and even save our data as CSV. It will allow us to easily get all images for the books.

How to use Scrapy?

- 1) Initialise a project
- 2) Initialise a spider with a start URL.
- 3) Write the scraper
- 4) Run the scraper

Spider in web scraping →

- 10 In web scraping, a spider is a function that allows us to visit all pages you want to scrape.

## Data Visualization

### Intro to DV

Creating charts and graphs for easy data exploration and understanding.

We can use Excel, Power BI, Tableau etc.

Advantage of doing data visualization in python are bigger amount of usable data, ability to integrate visualizations into an API or website.

25

How to do Data Visualization in Python?

Main library for data viz in Python is Matplotlib

Scatter Plot -

- 30 The most basic type of plot.
- Shows points in a 2D grid
- Useful to see
- Patterns and relationships
- See clusters and patterns.

### Parts of Plot

- Figure
- Axes (1 pair of X axis & Y axis are called as Axes).
- <sup>5</sup> X axis & Y axis
- Legend.

### Bar Plot

- 10 Plots for number vs category
- Bar plot is useful when visualising relationships between a category and income.  
eg Gender VS Income

### Histograms

Histograms are almost vertical bar plots where the height of the bar denotes the count of the category.

- 20 Histograms of numerical data can be made using Binning.

$$\text{bin size} = \frac{\max(x) - \min(x)}{\text{bins}}$$