#### Chapter 6:AI developement tools





Table of contents

- 1. Google Colab
- 2. Jupyter
- 3. Numpy
- 4. Matplotlib, SciKits Learn et Pandas.
- 5. Exercices
- 6.Conclusion

# Introduction:

- In this chapter, we will explore some of the most widely used tools and libraries in AI development. These tools are essential for data manipulation, visualization, and machine learning model development.
- Google Colab
- Jupyter
- > NumPy
- Matplotlib
- SciKit-Learn
- Pandas

#### Introduction:

Python programming concepts for data analysis, using popular libraries such as NumPy, Pandas, Matplotlib, and Scikit-learn. By the end of this course, you'll be comfortable with basic Python programming, handling data structures like NumPy arrays and Pandas DataFrames, and visualizing data with Matplotlib. Additionally, you'll gain an understanding of key tools used in AI development, including Google Colab and Jupyter Notebooks.

#### Google colab

1



#### Google colab :

- Google Colab (short for Colaboratory) is a free, cloudbased platform that allows you to write and execute Python code in a Jupyter notebook environment. It is particularly popular for AI and machine learning projects because it provides free access to GPUs and TPUs, making it easier to train complex models.
- ≻Key Features:
- Free access to GPUs and TPUs.
- Pre-installed libraries like TensorFlow, PyTorch, and Keras.

#### Google colab :

- Easy sharing and collaboration.
- Integration with Google Drive for saving and loading notebooks.
  Use Cases:
  - 1. Prototyping machine learning models.
  - 2. Running data analysis and visualization.
  - 3. Collaborating on AI projects with team members

2

# Jupyter



#### Jupyter :

Jupyter is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations, and narrative text. It supports multiple programming languages, but it is most commonly used with Python.

# Jupyter :

#### ≻Key Features:

- Interactive coding environment.
- Support for Markdown and LaTeX for documentation.
- Easy integration with data science libraries.
- Ability to export notebooks in various formats (HTML, PDF, etc.).

#### Jupyter :

≻Use Cases:

- Data cleaning and preprocessing.
- Exploratory data analysis (EDA).
- Prototyping and testing machine learning models.

2

# NumPy

# NumPy (Numerical Python)

NumPy (Numerical Python) is a fundamental library for scientific computing in Python. It provides support for arrays, matrices, and many mathematical functions to operate on these data structures.

- Efficient array operations.
- Broadcasting and vectorization for performance optimization.
- Integration with other libraries like SciPy, Pandas, and Matplotlib

- import numpy as np # import the library numpy
- array = np.array([1, 2, 3, 4, 5]) # Create a 1D array
- matrix = np.array([[1, 2], [3, 4]]) # Create a 2D array
- This imports the NumPy library and gives it the alias np. NumPy is a powerful library for numerical computations in Python.
- Why np? Using np as an alias is a convention in the Python community to make the code shorter and easier to read.

- print(array)
- output:
- [1 2 3 4 5]
- example :
- import numpy as np
- # Create an array of random temperatures in Celsius between -10 and 40
- temperatures = np.random.uniform(-10, 40, 5)
- print(temperatures)

- np.random.uniform(): Generates random values.
- np.array(): Creates an array.
- np.mean(), np.median(): Statistical operations.

#### MatPlotLib

4

# MatPlotLib:

- Matplotlib is a plotting library for Python that provides a wide variety of static, animated, and interactive visualizations. It is highly customizable and integrates well with other data science libraries.
- Support for a wide range of plot types (line, bar, scatter, histogram, etc.).
- Customizable plots with labels, legends, and annotations.
- Integration with Jupyter notebooks for online plotting

# MatPlotLib :

import matplotlib.pyplot as plt

# Sample data: Months and average temperatures months = ['January', 'February', 'March', 'April'] avg\_temperatures = [5, 7, 12, 16]

# Line plot of temperatures over the months

plt.plot(months, avg\_temperatures)

plt.xlabel('Months')

plt.ylabel('Average Temperature (°C)')

plt.title('Average Temperatures Over the Year')
plt.show()

## MatPlotLib :

```
Key Functions:

plt.plot(): Line plot.

plt.bar(), plt.scatter(): Bar chart and scatter plot.

plt.show(): Displays the plot.
```

- SciKit-Learn is one of the most popular machine learning libraries in Python, offering a wide range of algorithms for classification, regression, clustering, and more.
  - Key Features:
  - Supervised and unsupervised learning algorithms.
  - Tools for model evaluation and validation.
  - Easy integration with Pandas and NumPy.

• Example:

from sklearn.linear\_model import LinearRegression from sklearn.model\_selection import train\_test\_split

- # Dummy data for linear regression
- X = [[1], [2], [3], [4], [5]]
- y = [2, 4, 6, 8, 10]

# Split data into training and test sets

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2)

```
# Train a linear regression model
model = LinearRegression()
model.fit(X_train, y_train)
# Predict the test data
predictions = model.predict(X_test)
print(predictions)
```

- Key Functions:
- train\_test\_split(): Splits data into training and test sets.
- LinearRegression(): Linear regression model.
- model.fit(), model.predict(): Model training and prediction.

- Data Setup:
- We have the following data:
- Features (X): [[1], [2], [3], [4], [5]]
- Target (y): [2, 4, 6, 8, 10]
- This is a simple linear relationship where y is exactly double the value of X. So, for any given value of X, the value of y should be 2 \* X.

- Step 1: Splitting the Data
- Using train\_test\_split with test\_size=0.2, we randomly split the data into training and testing sets. Since the dataset is very small, there is only one data point in the test set and four in the training set. The actual split will be random, but here is an example of how the data might be split:

```
Training Data (X_train, y_train):

X_train = [[3], [5], [2], [4]]

y_train = [6, 10, 4, 8]

Test Data (X_test, y_test):

X_test = [[1]]

y_test = [2]
```

- Step 2: Training the Model
- The LinearRegression model will learn the relationship between the independent variable X (features) and the dependent variable y (target). It will compute the intercept and slope of the best-fit line.
- The expected values for the slope (β1) and intercept (β0) are:
- Slope ( $\beta$ 1): 2 (since y = 2 \* X in the dataset).
- Intercept (β0): 0 (since the line passes through the origin in this case).

- Step 3: Making Predictions
- After the model is trained, we can make predictions on the test set. In this case, X\_test = [[1]], and the model should predict a y value close to 2 (because y = 2 \* 1).
- Predictions: [2.]
- The model predicts y = 2.0 when X = 1, which is exactly what we expect because the relationship in the dataset is linear with a slope of 2 and an intercept of 0.

# Linear regression: example

- from sklearn.linear\_model import LinearRegression
- from sklearn.model\_selection import train\_test\_split
- import numpy as np
- # Example dataset
- experience = np.array([1, 2, 3, 4, 5]).reshape(-1, 1) # Feature
- salary = np.array([30000, 35000, 40000, 450000, 50000]) # Target

# Linear regression: example

- # Split data
- X\_train, X\_test, y\_train, y\_test = train\_test\_split(experience, salary, test\_size=0.2)
- # Create and train model
- model = LinearRegression()
- model.fit(X\_train, y\_train)
- # Predict
- predictions = model.predict(X\_test)

# Linear regression: example

- # Coefficient and Intercept
- print(f"Slope (Coefficient): {model.coef\_[0]}")
- print(f"Intercept: {model.intercept\_}")

# Pandas

#### Pandas :

- Pandas is an essential library for data manipulation and analysis, particularly with structured data. It provides DataFrame and Series objects for handling various types of datasets, from CSVs to SQL queries.
- Key Features:
- Handling and manipulation of structured data.
- Powerful operations such as groupby, pivot, and merge.
- Handling missing data and time series data.

#### Pandas :

import pandas as pd

# Create a DataFrame

data = {'Name': ['John', 'Anna', 'Peter'],

'Age': [28, 24, 35],

'City': ['New York', 'Paris', 'Berlin']}

df = pd.DataFrame(data)

# Display the DataFrame
print(df)

#### Pandas :

# Display descriptive statistics
print(df.describe())
Key Functions:
pd.DataFrame(): Creates a DataFrame.
df.groupby(): Group data by a specific column.
df.describe(): Summary statistics for numerical columns.

# Select a specific column ages = df['Age']# Filter rows where Age > 30 older\_than\_30 = df[df['Age'] > 30] # Create a new column  $df['Age_in_10_years'] = df['Age'] + 10$ # Sort the dataset by Age sorted\_df = df.sort\_values(by='Age')

- # Group by City and calculate the mean Age grouped = df.groupby('City')['Age'].mean()
- # Drop missing values
- df\_cleaned = df.dropna()
- # Fill missing values
- df\_filled = df.fillna({'Age': df['Age'].mean()})

- import pandas as pd
- import numpy as np
- # Create a DataFrame with missing values
- df = pd.DataFrame({
- 'Name': ['Alice', 'Bob', np.nan, 'David'],
- 'Age': [24, np.nan, 22, 32],
- 'City': ['New York', 'Paris', 'London', np.nan]})

- # Fill missing values
- df\_filled = df.fillna({'Name': 'Unknown', 'Age': df['Age'].mean(), 'City': 'Unknown'})
- print(df\_filled)

# Reading a Dataset with Pandas

import pandas as pd

# Read a dataset from a CSV file

- df = pd.read\_csv('sample\_data.csv')
- # Display the first 5 rows

print(df.head())

# Display basic information about the dataset
print(df.info())

# Display basic statistics
print(df.describe())

## Reading a Dataset with Pandas

read\_csv() reads the CSV file.

head() shows the first few entries.

info() gives a summary (columns, types, missing values).

describe() gives mean, std, min, max, etc

import pandas as pd

- # Reading a dataset from a CSV file
- # Example file: students\_scores.csv
- # Columns: 'Hours\_Studied', 'Test\_Score', 'Passed'
- df = pd.read\_csv('C:\\Users\\hp\\Desktop\\mémoire licence\\students\_scores.csv')
- # Display the first few rows
- print(df.head())
- # Display dataset information

print(df.info())
# Check for missing

# Check for missing values

print(df.isnull().sum())

from sklearn.model\_selection import train\_test\_split from sklearn.linear\_model import LinearRegression from sklearn.metrics import mean\_squared\_error import matplotlib.pyplot as plt

# Select features and target variable

- X = df[['Hours\_Studied']] # Independent variable
- y = df['Test\_Score'] # Dependent variable
- # Split into train and test sets
- X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2)
- # Create and train the model
- linear\_model = LinearRegression()
- linear\_model.fit(X\_train, y\_train)

# Make predictions

y\_pred = linear\_model.predict(X\_test)

# Evaluate the model

mse = mean\_squared\_error(y\_test, y\_pred)
print(f"Mean Squared Error (MSE): {mse:.2f}")
# Visualize the regression line
plt.scatter(X, y, color='blue')
plt.plot(X, linear\_model.predict(X), color='red')

plt.xlabel('Hours Studied')

plt.ylabel('Test Score')

plt.title('Linear Regression: Hours Studied vs Test Score')
plt.show()