# MACHINE LEARNING

**Instructor:** Ishaya, Jeremiah Ayock

**Lecture 03**: Notions and Definitions

January 20, 2025

Academic City University College, Agbogba Haatso, Ghana.

### Notations and Definitions

- Python
- Numpy
- Pandas
- Scipy
- matplotlib/seaborn
- Data Exploration

# Python

## PYTHON

Python is a widely used programming language for data science, machine learning, and artificial intelligence. It has become the go-to language for data scientists because of its simplicity, ease of use, and availability of various libraries and tools for data manipulation, visualization, and modelling.

In this lecture, we will cover some important Python libraries for data science, including **NumPy, pandas, SciPy, Matplotlib/Seaborn, and data exploration.**

## IMPORTANT PYTHON CONCEPTS FOR A DATA SCIENTIST:

- **Data scientists** need to be familiar with data structures such as **lists, tuples, dictionaries, and sets** to store, manipulate and organize data in Python.

- **Functions:**
  Functions are reusable blocks of code that perform a specific task. Data scientists need to be proficient in creating, defining and calling functions to reuse the same block of code in different parts of the program.

- **Control Flow:**
  Data scientists need to be familiar with control flow constructs like **loops and conditional statements** to execute the code in a specific sequence.

- **Modules and Libraries:**
  Python provides many libraries and modules to perform complex tasks. Data scientists should be familiar with libraries like NumPy, pandas, Matplotlib/Seaborn, SciPy, Scikit-Learn, and others.

- **File Input and Output:**
  Data scientists need to be familiar with file input/output (I/O) operations to read and write data to and from files.

- **Object-Oriented Programming (OOP):**
  Data scientists should have a basic understanding of OOP concepts like classes, objects, inheritance, and polymorphism. OOP is useful in creating custom data structures and functions for specific use cases.

- **Regular Expressions:**
  Regular expressions are useful in pattern matching and data cleaning operations. Data scientists should be familiar with the basic syntax and use cases of regular expressions.

- **Error Handling:** Data scientists should be proficient in error-handling techniques to debug their code and avoid crashes.

- **Virtual Environments:**
  Data scientists need to be familiar with virtual environments to create isolated environments for different projects with different dependencies.

- **Collaboration:**
  Data scientists should be familiar with version control systems like Git and collaboration platforms like GitHub to collaborate with other team members and manage the codebase.

# Numpy

**NumPy** is a **fundamental library** for scientific computing with Python.

It provides support for multidimensional arrays, as well as a wide range of mathematical functions to operate on them.

Here are some of the important NumPy concepts for a data scientist:

## NUMPY CON'T

- ndarrays NumPy's ndarrays are the core data structure for storing and manipulating numerical data. They provide efficient storage and access to homogeneous numerical data, with support for various mathematical operations. You should be familiar with **creating ndarrays, indexing and slicing them, and performing basic operations like addition, multiplication, and broadcasting.**

- Array Shape and Dimensions Understanding the shape and dimensions of an ndarray is essential for working with NumPy. The shape of an ndarray specifies the size of each dimension, while the number of dimensions is called the rank. You should be familiar with functions for manipulating the shape and dimensions of ndarrays, like **reshape(), resize(), and transpose().**

## NUMPY CON'T

- Array Indexing and Slicing
  Indexing and slicing are essential operations for accessing and modifying elements of an array. You should be familiar with basic indexing and slicing, as well as advanced indexing techniques like Boolean indexing, fancy indexing, and broadcasting.

- Array Operations NumPy provides a wide range of mathematical and logical operations to work with ndarrays. You should be familiar with functions like **np.sum(), np.mean(), np.max(), np.min(), np.std(), np.dot(), np.transpose(), np.concatenate()**, and many more.

- **Broadcasting**
  Broadcasting is a powerful feature of NumPy that allows for element-wise operations between ndarrays of different shapes and dimensions.

- **Vectorization**
  Vectorization is the process of converting iterative operations into vector operations, which can be performed more efficiently with NumPy.

- **File Input/Output**
  NumPy provides functions for reading and writing ndarrays to and from files. You should be familiar with functions like **np.load(), np.save(), np.savetxt(), and np.loadtxt()** for working with files in NumPy.

# pandas

**Pandas** is a popular Python library for data manipulation and analysis. It provides high-performance data structures for efficiently working with large datasets, as well as a wide range of functions for data cleaning, exploration, and transformation. Here are some of the important panda's concepts for a data scientist:

## PANDAS CON'T

- **Series and DataFrame**

  Pandas provide two main data structures for working with tabular data: **Series and DataFrame**. A Series is a one-dimensional labeled array that can hold any data type, while a data frame is a two-dimensional labeled data structure with columns of potentially different data types.

- **Indexing and Selection**

  Pandas provides various indexing and selection methods for accessing and modifying elements of a Series or DataFrame. You should be familiar with basic indexing and selection methods like **loc, iloc, and ix,** as well as advanced techniques like Boolean indexing, fancy indexing, and hierarchical indexing.

- **Data Cleaning and Preprocessing** Data cleaning and preprocessing are essential steps in data analysis. Pandas provides a wide range of functions for handling missing data, removing duplicates, renaming columns, converting data types, and more. You should be familiar with functions like **dropna(), fillna(), drop_duplicates(), rename(), astype(), and apply()** for cleaning and preprocessing data.

- **Data Exploration and Analysis**
  Pandas provide a wide range of functions for exploring and analyzing data. You should be familiar with functions for calculating summary statistics like mean, median, and standard deviation, as well as functions for aggregating, grouping, and pivoting data. You should also be familiar with functions for working with time-series data, like **resample() and rolling()**.

- **Data Transformation**
  Data transformation involves converting data from one
  form to another. Pandas provide a wide range of
  functions for data transformation, including functions for
  sorting, ranking, merging, and pivoting data. You should
  be familiar with functions like **sort_values(), rank(),
  merge(), and pivot_table()** for transforming data.

## PANDAS CON'T

- **Time-Series Analysis**
  Pandas provide powerful support for time-series analysis.
  You should be familiar with functions for working with
  time-series data, like **resample(), rolling(), and shift()**.
  You should also be familiar with functions for handling
  time zones and date ranges.

- **Input/Output**
  Pandas provide functions for reading and writing data to
  and from various file formats, including **CSV, Excel,
  SQL databases,** and more. You should be familiar with
  functions like **read_csv(), read_excel(), read_sql(),
  to_csv(), and to_excel()** for working with data in
  Pandas.

# Scipy

**Scipy** is a powerful Python library for scientific computing that provides many functions for **numerical optimization, integration, interpolation, signal processing, linear algebra, and more.** Here are some of the important Scipy concepts for a data scientist:

## SCIPY

- **Integration**
  Scipy provides functions for numerical integration, including **quad(), dblquad(), and tplquad()**. These functions can be used to calculate integrals of functions in one, two, or three dimensions, respectively.

- **Optimization**
  Scipy provides functions for numerical optimization, including **minimize(), curve_fit(), and root()**. These functions can be used to find the minimum or maximum of a function, fit a curve to data, or solve nonlinear equations, respectively.

## SCIPY

- **Interpolation**
  Scipy provides functions for numerical interpolation, including **interp1d(), interp2d(), and griddata()**. These functions can be used to interpolate data onto a grid, or to create a smooth curve that passes through a set of points.

- **Signal Processing**
  Scipy provides functions for signal processing, including **convolution(), fft(), and spectrogram()**. These functions can be used to **filter, transform, and analyze signals, such as audio, image, or time-series** data.

## SCIPY

- **Linear Algebra**
  Scipy provides functions for linear algebra, including **solve(), eig(), and svd()**. These functions can be used to solve linear systems of equations, compute eigenvalues and eigenvectors, and perform singular value decomposition.

- **Statistics**
  Scipy provides functions for statistical analysis, including **ttest_1samp(), ttest_ind(), and pearsonr()**. These functions can be used to perform hypothesis testing, calculate confidence intervals, and compute correlation coefficients.

## SCIPY

- **Sparse Matrices**
  Scipy provides support for sparse matrices, which are useful for representing large datasets with many zeros. Scipy provides functions for creating, manipulating, and solving sparse matrices, including **csr_matrix(), coo_matrix(), and spsolve()**.

- **Image Processing**
  Scipy provides functions for image processing, including **imread(), imsave(), and ndimage()**. These functions can be used to read and write image files, as well as perform operations like filtering, segmentation, and morphological operations on images.

# matplotlib/seaborn

**Matplotlib** is a popular Python library for creating visualizations and plots. Here are some of the important Matplotlib concepts for a data scientist:

# MATPLOTLIB CON'T

- **Figure and Axes Objects**
  Matplotlib operates on two main types of objects: **Figure and Axes**. The Figure object represents the entire figure or window, while the Axes object represents an individual plot within the figure. Understanding how to create, customize, and manipulate these objects is key to creating effective visualizations.

- **Types of Plots**
  Matplotlib supports a wide range of plot types, including line plots, scatter plots, bar plots, histogram plots, and more. Understanding the syntax and options for each plot type is important for creating the desired visualizations.

## MATPLOTLIB CON'T

- **Subplots**
  Matplotlib allows you to create multiple plots within a
  single figure using subplots. Understanding how to create
  and customize subplots can be useful for comparing
  multiple datasets or visualizing different aspects of a
  single dataset.

- **Saving and Exporting Plots**
  Matplotlib allows you to save your plots in various
  formats, such as **PNG, PDF, or SVG**. Understanding
  how to save and export your plots can be useful for
  sharing your visualizations with others or incorporating
  them into reports or presentations.

- **Plot Customization**
  Matplotlib provides many options for customizing plots, such as changing the color, size, and style of lines or markers, adding labels and titles, adjusting axis limits and ticks, and more. Understanding how to use these options can help improve the clarity and effectiveness of your visualizations.

- Integration with Pandas Matplotlib can be easily integrated with the Pandas library, which is commonly used for data manipulation and analysis. Understanding how to use Matplotlib to create visualizations from Pandas dataframes can be useful for quickly exploring and analyzing datasets.

# Data Exploration

Data exploration is a critical step in the data science process, as it involves understanding the structure, quality, and patterns in the data. Here are some of the important data exploration concepts for a data scientist:

## DATA EXPLORATION

- **Data Types and Formats**
  Understanding the types and formats of data is important for determining appropriate analysis methods and identifying potential data quality issues. Common data types include numerical, categorical, text, and datetime data.

- **Descriptive Statistics**
  Descriptive statistics provide summary information about the data, such as measures of central tendency (e.g., mean, median) and variability (e.g., standard deviation, range). Understanding how to calculate and interpret descriptive statistics is important for identifying patterns and outliers in the data.

- **Data Visualization**

## DATA EXPLORATION

- **Data Cleaning and Preprocessing**
  Data cleaning involves identifying and correcting errors, missing values, and outliers in the data, while data preprocessing involves transforming and normalizing the data to prepare it for analysis. Understanding how to use tools like pandas and numpy for data cleaning and preprocessing is critical for ensuring the quality and reliability of analysis results.

- **Correlation Analysis**
  Correlation analysis involves identifying the strength and direction of the relationships between variables in the dataset. Understanding how to calculate and interpret correlation coefficients can help identify important variables and relationships within the data.

## DATA EXPLORATION

- **Feature Engineering**
  Feature engineering involves creating new features or
  variables from existing data to improve the performance
  of machine learning models. Understanding how to select
  and create appropriate features is important for
  developing effective models.

- **Exploratory Data Analysis (EDA)** EDA involves
  examining the data in depth to generate hypotheses and
  insights about the data. Techniques such as clustering
  and dimensionality reduction can help identify patterns
  and relationships in the data.

- **Data Quality Assessment**
  Assessing the quality of the data is critical for ensuring the reliability and validity of analysis results. Common methods for assessing data quality include evaluating data completeness, accuracy, and consistency.

Overall, effective data exploration is critical for understanding the characteristics of the data and identifying potential issues or patterns that can inform subsequent analysis steps.

# END OF PRESENTATION

THANK YOU