# Notes on: Python libraries suitable for Machine Learning, Numpy_from_0

## 1.) Introduction

Introduction to Machine Learning (ML) is about enabling computers to learn from data to make predictions or decisions without explicit programming. For example, ML powers recommendation systems, facial recognition, and medical diagnostics. Effectively implementing these techniques relies heavily on powerful programming languages and specialized libraries. Python has become the industry standard for ML due to its straightforward syntax, readability, and a vast collection of libraries that simplify complex data science and ML tasks.

1- Python's Dominance in Machine Learning
Python's widespread adoption in ML is due to several factors:
   • Its clear and concise syntax allows developers to write less code for complex operations.
   • A large, supportive community contributes to extensive documentation and resources.
   • The availability of a rich ecosystem of high-quality, open-source libraries optimized for numerical computing, data analysis, and machine learning algorithms.

2- Python Libraries for Machine Learning
The Python ML ecosystem is built upon a stack of inter-dependent libraries. Understanding these foundational libraries is crucial. One of the most fundamental and frequently used libraries is NumPy. It serves as a bedrock for almost all other scientific and numerical computing libraries in Python, including those directly involved in machine learning.

3- NumPy: The Core of Numerical Computing
NumPy, short for Numerical Python, is a library designed to provide high-performance operations on large, multi-dimensional arrays and matrices. It also offers a comprehensive collection of mathematical functions to efficiently manipulate these arrays.

   • What it is for: NumPy is optimized to handle numerical data much faster and with greater memory efficiency than standard Python lists. This capability is paramount when working with the massive datasets typical in machine learning.
   • The central concept: At the heart of NumPy is the 'ndarray' (N-dimensional array) object. An ndarray is a grid of values, all of the same data type, indexed by a tuple of non-negative integers. It can represent anything from a simple list of numbers (1D array) to a spreadsheet-like structure (2D array/matrix), or even higher dimensions for complex data.
   • Advantages over Python lists:
   • Performance: NumPy operations are implemented in optimized C code, making them significantly faster for numerical computations compared to equivalent operations on standard Python lists.
   • Memory Efficiency: Ndarrays store elements of the same type contiguously in memory, leading to lower memory consumption, especially beneficial for large datasets.
   • Rich Functionality: NumPy provides powerful tools for linear algebra, Fourier transforms, random number generation, and other mathematical operations essential for building and optimizing ML algorithms.

4- Real-World Impact of NumPy in ML
Consider an ML task like processing sensor data from a smart device. This data might come as a continuous stream of numerical readings, forming a time series. Representing this data as a NumPy ndarray allows for highly efficient computations—like calculating averages, finding deviations, or performing transformations—across millions of data points quickly. Without NumPy, these fundamental data manipulations would be prohibitively slow, making complex ML models impractical. It provides the raw computational power needed to process the numerical backbone of most machine learning models.

5- NumPy as a Foundational Layer

Many other critical Python libraries used in ML, such as Pandas (for structured data analysis), Scikit-learn (for various machine learning algorithms), and even deep learning frameworks like TensorFlow and PyTorch, either build upon NumPy or are designed to seamlessly integrate with NumPy arrays. Therefore, a solid grasp of NumPy is a prerequisite for advancing in machine learning with Python.

Summary of Key Points:
- Machine Learning uses data to learn and make predictions.
- Python is a leading language for ML due to its powerful library ecosystem.
- NumPy is a foundational Python library for efficient numerical computing.
- Its core object is the 'ndarray', which stores multi-dimensional, same-type numerical data.
- NumPy offers significant performance and memory advantages over Python lists for numerical operations, crucial for large datasets.
- It enables efficient representation and manipulation of numerical data that underpins most ML algorithms, like image pixels or sensor readings.
- NumPy is a fundamental building block for many other advanced ML and scientific computing libraries.

# 2.) Creating Array: array()

Creating Array: array() in NumPy

NumPy is a cornerstone Python library in Machine Learning, providing highly efficient array objects that are essential for handling numerical data. While Python lists can store collections of data, NumPy's N-dimensional array (ndarray) is specifically designed for high-performance numerical operations crucial for ML algorithms.

1. What is a NumPy Array?
- A NumPy array is a grid of values, all of the same data type, indexed by a tuple of non-negative integers. It is a fundamental data structure in numerical computing.
- Unlike Python lists, NumPy arrays are homogeneous, meaning all elements must be of the same data type (e.g., all integers, all floats). This homogeneity allows for highly optimized operations.

2. The array() Function
- The primary way to create a NumPy array from existing Python data structures, such as lists or tuples, is by using the numpy.array() function.
- Its main purpose is to convert an input (like a list) into an ndarray, allowing you to leverage NumPy's powerful features.

3. Creating Arrays from Python Lists
- To use array(), you first need to import NumPy, usually aliased as np.
- Example: import numpy as np

a. One-Dimensional Arrays (Vectors)
- These are like a single row or column of numbers, often representing a list of features for a single sample in ML.
- Example: my_list = [10, 20, 30, 40]
- vector = np.array(my_list)
- Result: [10 20 30 40]

b. Two-Dimensional Arrays (Matrices)
- Matrices are grids with rows and columns, commonly used to represent datasets in ML, where rows are samples and columns are features.
- You pass a list of lists to array().
- Example: data = [[1, 2, 3], [4, 5, 6]]
- matrix = np.array(data)

• Result:
[[1 2 3]
[4 5 6]]

c. Higher-Dimensional Arrays (Tensors)
    • NumPy can create arrays with three or more dimensions. In ML, 3D arrays are often used for image data (height, width, color channels) or sequences.
    • Example: image_data = [[[0,0,0], [1,1,1]], [[2,2,2], [3,3,3]]] (a simplified 2x2 pixel image with 3 color channels)
    • tensor = np.array(image_data)
    • This creates a 3D array, which is conceptually a **tensor** in ML.

4. Understanding Data Types (dtype)
    • When you create an array, NumPy infers the most suitable data type for its elements. For instance, if you provide integers, it might choose int32 or int64. If you mix integers and floats, it will infer float64 to avoid data loss.
    • You can explicitly specify the data type using the dtype argument. This is crucial for memory management and numerical precision in ML applications.
    • Example: float_array = np.array([1, 2, 3], dtype=float) or dtype='float64'
    • Example: int_array = np.array([1.5, 2.7], dtype=int) (note: this truncates decimal parts)
    • Common dtypes: int, float, bool, str. Specific types like int8, float32 offer finer control.

5. Essential Array Attributes (upon creation)
    • After creating an array, it has important attributes that describe its structure.
    • .shape: A tuple indicating the size of the array in each dimension.
    • .ndim: The number of dimensions (axes) in the array.
    • .size: The total number of elements in the array.
    • These attributes help you understand the array's structure, which is fundamental for preparing data for ML models.

6. Real-World Relevance in Machine Learning
    • Input Features: Datasets for ML models are typically converted into 2D NumPy arrays (matrices) where each row is a sample and each column is a feature.
    • Image Processing: Images are often represented as 3D (height, width, color channels) or 4D (number of images, height, width, channels) NumPy arrays.
    • Model Weights: The parameters of machine learning models (like weights in neural networks) are stored and manipulated as NumPy arrays.

Summary of Key Points:
    • np.array() is the fundamental function to create NumPy arrays from Python lists or tuples.
    • NumPy arrays are homogeneous, meaning all elements share the same data type.
    • You can create 1D (vectors), 2D (matrices), and higher-dimensional arrays (tensors).
    • The dtype argument allows explicit control over the array's data type, important for precision and memory.
    • Attributes like shape, ndim, and size describe the array's structure immediately after creation, which is vital for ML data handling.

# 3.) Accessing Array: by referring to its index number

Accessing elements in a Numpy array by referring to their index number is a fundamental operation. Once data is stored in an array, we often need to retrieve, modify, or analyze specific pieces of information. This process is called indexing.

1. Understanding Array Indexing

    • Arrays organize data in a sequential or structured manner.

• Each element within an array is assigned a unique position number, known as an index.
• In Python, and consequently in Numpy, indexing is '0-based'. This means the first element is at index 0, the second at index 1, and so on.
• Think of it like house numbers on a street, starting from zero: house number 0, house number 1, etc.

## 2. Accessing Elements in One-Dimensional (1D) Arrays

• For a 1D array (a vector), you access elements using a single index within square brackets.
• Syntax: array_name[index]
• Example:

```
import numpy as np
data_points = np.array([10, 20, 30, 40, 50])
first_element = data_points[0] # Accesses 10
third_element = data_points[2] # Accesses 30
```

• This allows you to pinpoint and retrieve a single value from the sequence.

## 3. Accessing Elements in Multi-Dimensional Arrays

• For multi-dimensional arrays (like 2D matrices or 3D tensors), you need an index for each dimension.
• For a 2D array (matrix), you specify the row index followed by the column index.
• Syntax: array_name[row_index, column_index]
• Example (2D array):

```
matrix = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
element_at_row1_col2 = matrix[1, 2] # Accesses 6 (row 1, column 2)
# Think of a spreadsheet: [row number, column number] to get a specific cell.
```

• For 3D arrays, you'd use three indices: array_name[dim1_index, dim2_index, dim3_index]. This is crucial for handling complex data like images (height, width, color channels) in ML.

## 4. Negative Indexing

• Numpy also supports negative indexing, which allows you to access elements from the end of the array.
• The last element is at index -1, the second to last at -2, and so on.
• Syntax: array_name[-index_from_end]
• Example:

```
data_points = np.array([10, 20, 30, 40, 50])
last_element = data_points[-1] # Accesses 50
second_to_last = data_points[-2] # Accesses 40
```

• This is particularly useful when you need to quickly grab the most recent data point or a target label located at the end of a feature set.

## 5. Importance in Machine Learning

• Accessing specific elements is fundamental for preparing and manipulating data for ML models.
• It allows you to extract individual features, retrieve a specific data sample, or isolate target labels from a dataset.
• For instance, if your dataset array contains features in the first columns and the target variable in the last, you'd use indexing to separate them for training.
• It is also critical for inspecting model predictions, error analysis, and updating specific parameters.

Summary of Key Points:
• Indexing is how you retrieve specific elements from an array.
• Numpy uses 0-based indexing (first element is at index 0).
• Use array_name[index] for 1D arrays.
• Use array_name[row_index, column_index] for 2D arrays, extending for higher dimensions.
• Negative indexing (e.g., array_name[-1]) accesses elements from the end of the array.
• This skill is vital for data extraction, manipulation, and analysis in Machine Learning workflows.

# 4.) Stacking & Splitting: stack(), array_split()

In the context of Machine Learning, data often needs to be organized or segmented in specific ways. Numpy's **stacking** and **splitting** functions provide powerful tools for these operations, essential for tasks like preparing input data, combining different features, or dividing datasets for training.

1. Importance of Stacking and Splitting in Machine Learning

   • Data Preparation: Machine Learning models expect data in specific shapes. Stacking helps combine individual data points or features, while splitting helps divide large datasets into manageable chunks.
   • Feature Engineering: Stacking can be used to combine different types of features (e.g., sensor readings, image channels) into a single, comprehensive data representation for a model.
   • Training Efficiency: Splitting datasets allows for creating training, validation, and test sets. It also enables techniques like mini-batch gradient descent, where a large dataset is processed in smaller batches.

2. Stacking Arrays with numpy.stack()

   • What it is: The numpy.stack() function joins a sequence of arrays along a new axis. Unlike concatenate, which joins along an existing axis, stack adds a dimension.
   • How it works: All input arrays must have the same shape. When stacked, they are arranged along a new axis specified by the axis parameter. If axis=0 (default), the new dimension is added at the beginning; if axis=1, it's added as the second dimension, and so on.
   • Example: Imagine you have two separate lists of daily temperatures and humidity readings, each for 5 days.
   • import numpy as np
   • temps = np.array([20, 21, 22, 23, 24]) -- Shape (5,)
   • humidity = np.array([60, 62, 65, 63, 61]) -- Shape (5,)
   • weather_data = np.stack((temps, humidity), axis=0)
   • -- Resulting 'weather_data' will have shape (2, 5).
   • -- First row: temperatures, Second row: humidity.
   • -- If axis=1 was used, shape would be (5, 2), where each row is [temp, humidity].
   • Real-world relevance: Combining multiple sensor readings (e.g., temperature, pressure, light) into a single feature vector for a time step. Creating batches of images where each image is an individual array, and stacking them creates a batch with a new 'batch' dimension.

3. Splitting Arrays with numpy.array_split()

   • What it is: The numpy.array_split() function divides an array into multiple sub-arrays along a specified axis. It is more flexible than numpy.split(), as it allows for uneven divisions.
   • How it works: You provide the array to split and either the number of sub-arrays you want or a list of indices at which to split. array_split() tries to distribute the elements as evenly as possible if the number of sub-arrays does not divide the total number of elements perfectly.
   • Example: Dividing a dataset of 10 samples into 3 parts.
   • import numpy as np
   • full_dataset = np.arange(10).reshape(10, 1) -- A 10x1 array
   • train_val_test = np.array_split(full_dataset, 3)
   • -- 'train_val_test' will be a list of 3 arrays.
   • -- The first two arrays will have 4 rows each, and the last one will have 2 rows.
   • -- This demonstrates uneven splitting.
   • Real-world relevance: Dividing a large dataset into training, validation, and test sets for model evaluation. Creating mini-batches of data for iterative optimization algorithms like gradient descent, which process small subsets of data at a time. Distributing data across multiple processors for parallel computation.

4. Key Takeaways

   • numpy.stack() is for combining arrays along a *new* axis, useful for creating a consolidated view of related data or batches.
   • numpy.array_split() is for dividing a single array into multiple smaller ones along an *existing* axis, especially useful for dataset management and processing.
   • Both functions are fundamental for preparing and manipulating data efficiently for various Machine Learning tasks.


# 5.) Maths Functions: add(), subtract(), multiply(), divide(), power(), mod()

NumPy, a core Python library for numerical operations, is vital for Machine Learning due to its efficient handling of large arrays. It provides a suite of mathematical functions known as **Universal Functions** or ufuncs. These functions operate element-wise on arrays, significantly speeding up computations compared to standard Python loops, which is crucial when dealing with massive datasets in Machine Learning.

Here, we will explore some fundamental mathematical ufuncs:

1. numpy.add()
   • Explanation: Performs element-wise addition between two arrays or an array and a scalar. If adding two arrays, they must have compatible shapes.
   • Example: If you have an array A = [1, 2, 3] and array B = [4, 5, 6], numpy.add(A, B) results in [5, 7, 9].
   • ML Context: Used for combining different feature vectors, adjusting model parameters (like adding a bias term), or updating weights in algorithms during training.

2. numpy.subtract()
   • Explanation: Performs element-wise subtraction between two arrays or an array and a scalar.
   • Example: If A = [10, 20, 30] and B = [1, 2, 3], numpy.subtract(A, B) results in [9, 18, 27].
   • ML Context: Essential for calculating the **error** or **difference** between a model's prediction and the actual target value, a fundamental step in many optimization algorithms.

3. numpy.multiply()
   • Explanation: Performs element-wise multiplication between two arrays or an array and a scalar. This is not matrix multiplication.
   • Example: If A = [1, 2, 3] and B = [2, 3, 4], numpy.multiply(A, B) results in [2, 6, 12].
   • ML Context: Used for scaling features, applying weights to inputs in a neural network layer, or computing products in cost functions.

4. numpy.divide()
   • Explanation: Performs element-wise division between two arrays or an array and a scalar.
   • Example: If A = [10, 20, 30] and B = [2, 5, 10], numpy.divide(A, B) results in [5., 4., 3.].
   • ML Context: Often used for normalizing data (e.g., dividing by the maximum value or standard deviation) or calculating averages across certain dimensions.

5. numpy.power()
   • Explanation: Raises each element in the first array to the power of the corresponding element in the second array, or to a scalar power.
   • Example: If A = [2, 3, 4] and you want to square each element, numpy.power(A, 2) results in [4, 9, 16].
   • ML Context: Useful in feature engineering (creating polynomial features), calculating certain parts of cost functions (e.g., squared error), or activation functions.

6. numpy.mod() (also numpy.remainder())

• Explanation: Computes the element-wise remainder of the division of the first array by the second array (or a scalar).
    • Example: If A = [10, 11, 12] and B = 3, numpy.mod(A, B) results in [1, 2, 0].
    • ML Context: Less directly applied in core algorithms but useful in data pre-processing, for instance, when dealing with cyclical features (e.g., time of day) or data binning.

These ufuncs are optimized C implementations, making them significantly faster than Python's built-in operations on lists for large numerical data, which is fundamental to the performance of Machine Learning models.

Summary of Key Points:
    • NumPy ufuncs perform efficient element-wise mathematical operations on arrays.
    • add(), subtract(), multiply(), divide() handle basic arithmetic.
    • power() is used for exponentiation.
    • mod() calculates remainders.
    • These functions are crucial for data manipulation, feature engineering, and core calculations in Machine Learning algorithms, directly impacting performance.

# 6.) Statistics Functions: amin(), amax(), mean(), median(), std(), var(), average(), ptp()

Statistics Functions in Numpy for Machine Learning

Numpy is a fundamental library in Machine Learning (ML) primarily because of its efficient array operations. When dealing with large datasets, understanding the characteristics of your data is crucial. Statistical functions help us summarize and interpret these datasets, which are often stored as Numpy arrays. These functions are vital for data exploration, preprocessing, and understanding model performance.

Let's explore some key statistical functions provided by Numpy:

    • amin() - Minimum Value
This function finds the smallest value within a given Numpy array.
It helps identify the lower bound of your data, like the lowest recorded temperature or the minimum feature value.
Example: In a dataset of house prices, amin() would show the cheapest house.

    • amax() - Maximum Value
This function finds the largest value within a given Numpy array.
It helps identify the upper bound of your data, like the highest exam score or the maximum sensor reading.
Example: In a dataset of patient ages, amax() would show the oldest patient.

    • mean() - Arithmetic Mean (Average)
This calculates the arithmetic average of all values in an array. It sums all values and divides by the count of values.
The mean is a common measure of central tendency, indicating the typical value.
Example: Calculating the average accuracy of a machine learning model over several test runs.

    • median() - Middle Value
This finds the middle value in a sorted dataset. If the number of data points is even, it's the average of the two middle values.
The median is robust to outliers, meaning extreme values don't heavily influence it, making it useful for skewed data.
Example: When analyzing incomes, the median income gives a better sense of a typical income than the mean, as a few very high earners can skew the mean.

• std() - Standard Deviation
This measures the dispersion or spread of data points around the mean. A low standard deviation means data points are generally close to the mean, while a high standard deviation indicates wider spread.
It quantifies how much individual data points deviate from the average.
Example: In quality control, a low standard deviation of product weights indicates consistency in manufacturing.

• var() - Variance
This is the average of the squared differences from the mean. It's simply the square of the standard deviation.
Variance provides a measure of how far each number in the set is from the mean. It's often used in theoretical statistics and certain ML algorithms.
Example: Comparing the variance of stock prices can help assess risk; higher variance means higher volatility.

• average() - Weighted Average
This function calculates a weighted average, allowing certain data points to contribute more or less to the final average. You provide a 'weights' array alongside your data.
It's useful when different data points have different levels of importance or reliability.
Example: Calculating a student's GPA where different courses have different credit weights.

• ptp() - Peak-to-Peak (Range)
This computes the range of values, which is the difference between the maximum and minimum values in the array (amax - amin).
It gives a quick measure of the spread or variability of the entire dataset.
Example: Determining the fluctuation range of a stock price over a period.

Summary of Key Points:
  • Numpy's statistical functions are essential for understanding and preprocessing data in Machine Learning.
  • amin() and amax() provide the bounds of your data.
  • mean() and median() describe central tendency, with median being more robust to outliers.
  • std() and var() quantify the spread or variability of data around the mean.
  • average() offers a flexible weighted mean calculation.
  • ptp() quickly shows the total range of your data.


# 7.) Summary And Revision

Summary And Revision of Numpy for Machine Learning

1- Understanding the importance of revisiting foundational topics like Numpy is crucial for any computer engineering student diving into Machine Learning. This session aims to consolidate your understanding of Numpy beyond basic operations, focusing on its core strengths and how it prepares data for powerful ML algorithms.

2- Key Concepts for Deeper Understanding in Numpy:

2.1- Numpy's Efficiency and Vectorization
  • While you've learned basic operations, it's vital to grasp *why* Numpy is preferred. It's written in C, making operations incredibly fast compared to Python lists.
  • Vectorization is the core principle: applying operations to entire arrays at once, rather than element-by-element loops. This dramatically speeds up computations, which is essential for large datasets in ML.
  • This efficiency directly translates to faster training times for your ML models.

2.2- Array Attributes: Shape, Data Type (dtype), and Dimensions (ndim)
   • Beyond just creating and accessing arrays, understanding their properties is key.
   • The shape attribute (array.shape) tells you the size of the array along each dimension, like (rows, columns). This defines how data is structured.
   • The dtype attribute (array.dtype) specifies the data type of elements in the array (e.g., int32, float64). Consistent data types allow Numpy to store and process data efficiently.
   • The ndim attribute (array.ndim) indicates the number of dimensions (axes) an array has. A vector is 1-D, a matrix is 2-D, etc. These attributes are critical for debugging and ensuring data is in the correct format for ML libraries.

2.3- Reshaping and Transposing Data for ML
   • Machine Learning models often require data in specific shapes.
   • reshape() allows you to change the dimensions of an array without changing its data. For example, converting a 1D array into a 2D column vector, or vice versa, is common for model inputs.
   • Transposing (array.T) swaps rows and columns, which is essential in linear algebra and for preparing features or target variables in ML.
   • Example: A single image might need to be flattened into a 1D array, or a sequence of features reshaped to fit a recurrent neural network input.

2.4- Broadcasting: Simplifying Array Operations
   • Broadcasting is a powerful mechanism that allows Numpy to perform operations on arrays of different shapes.
   • When two arrays have compatible shapes, Numpy **stretches** the smaller array to match the larger one for the operation, without actually duplicating data in memory.
   • Example: Adding a scalar value to an array, or adding a 1D array (vector) to a 2D array (matrix) row-wise or column-wise.
   • This avoids writing explicit loops and makes code cleaner and more efficient.

2.5- Boolean Indexing and Filtering
   • This is a powerful way to select subsets of data based on conditions.
   • You create a boolean array (True/False) by applying a condition to your Numpy array.
   • This boolean array can then be used as an index to select only the elements where the condition is True.
   • Example: Filtering out all data points above a certain threshold, or selecting only positive values. This is invaluable for data cleaning and feature selection in ML.

3- Real-world application of these concepts involves data preprocessing, feature engineering, and preparing the exact input shape required by various ML algorithms. Mastering these advanced Numpy ideas ensures you can efficiently manipulate and prepare complex datasets for machine learning models.

4- Summary of Key Points:
   • Numpy's speed comes from C implementation and vectorization.
   • Understand shape, dtype, ndim for data structure and type.
   • reshape() and T are crucial for data transformation for ML models.
   • Broadcasting enables efficient operations on arrays of different shapes.
   • Boolean indexing allows powerful conditional data selection.