# Lesson 16: Performing Data Manipulation and Analysis using Python Libraries

When it comes to data analytics in Python, several key libraries provide powerful tools and functionalities. Among these, NumPy, pandas, and Matplotlib are widely used and form the foundation of many data analysis workflows. Let's explore these libraries and their roles in data analytics:

## 1. NumPy:

NumPy, short for Numerical Python, is a fundamental library for scientific computing in Python. It provides support for efficient numerical operations on large multidimensional arrays and matrices. NumPy's main data structure is the ndarray (n-dimensional array), which allows for fast and vectorized computations. Key features of NumPy include:

- Mathematical functions: NumPy offers a wide range of mathematical functions, such as trigonometric, exponential, and statistical functions, enabling complex numerical calculations.
- Array operations: NumPy provides extensive capabilities for array manipulation, including indexing, slicing, reshaping, and merging arrays.
- Broadcasting: NumPy's broadcasting feature allows for performing operations on arrays of different shapes and sizes, providing flexibility in computations.
- Linear algebra: NumPy includes linear algebra functions for matrix operations, such as matrix multiplication, decomposition, and solving linear equations.
- Integration with other libraries: NumPy seamlessly integrates with other data analysis libraries like pandas and Matplotlib, enabling smooth data processing and visualization.

## 2. pandas:

pandas is a powerful library built on top of NumPy that provides data manipulation and analysis tools. It offers high-performance data structures and data analysis functions, making it a go-to library for working with structured data. Key features of pandas include:

- DataFrame: The core data structure in pandas is the DataFrame, which represents tabular data with labeled rows and columns. It allows for efficient handling and manipulation of structured data.
- Data cleaning and preprocessing: pandas provides functions for cleaning, filtering, and transforming data, handling missing values, and performing data imputation.
- Data aggregation and grouping: pandas supports grouping data based on specific criteria and performing aggregations, such as sum, mean, count, etc., on grouped data.
- Time series analysis: pandas offers functionality for working with time series data, including date/time indexing, resampling, and time-based calculations.
- Data input/output: pandas supports reading and writing data from various file formats, such as CSV, Excel, SQL databases, and more.
- Integration with other libraries: pandas integrates well with other libraries, making it a crucial component in data analysis workflows alongside NumPy and Matplotlib.

# 3. Matplotlib:

Matplotlib is a powerful plotting library for creating static, animated, and interactive visualizations in Python. It provides a wide range of plotting functions and customization options to represent data effectively. Key features of Matplotlib include:

- Variety of plots: Matplotlib supports a comprehensive range of plot types, including line plots, scatter plots, bar plots, histograms, pie charts, and more.
- Customization: Matplotlib allows fine-grained control over plot aesthetics, such as colors, markers, line styles, titles, labels, and annotations.
- Subplots and layouts: Matplotlib enables the creation of multiple subplots within a single figure, allowing for side-by-side visualizations and complex layouts.
- Export and integration: Matplotlib supports exporting plots to various file formats, including PNG, PDF, SVG, etc. It can also be integrated with GUI frameworks like PyQt and web frameworks like Django and Flask.
- Seaborn integration: Seaborn, built on top of Matplotlib, provides a higher-level interface for creating visually appealing statistical visualizations with fewer lines of code.

These three libraries, NumPy, pandas, and Matplotlib, form a powerful trio for data analytics in Python. NumPy handles efficient numerical computations, pandas provides powerful data manipulation capabilities, and Matplotlib enables the creation of high-quality visualizations. Together, they offer a comprehensive toolkit for data

analysis, from data cleaning and preprocessing to exploratory data analysis and visualization.

# Working with NumPy arrays

## 1. Creating NumPy Arrays:

NumPy arrays can be created in various ways:

- ***From a list or tuple using the `numpy.array()` function:***

You can convert a list or tuple into a NumPy array by passing it as an argument to **`np.array()`**.

```python
import numpy as np

my_list = [1, 2, 3, 4, 5]
arr = np.array(my_list)
print(arr)   # Output: [1 2 3 4 5]
```

- ***Using built-in functions like `numpy.zeros()` or `numpy.ones()`:***

You can create arrays of specific sizes with zeros or ones using these functions.

```python
zeros_arr = np.zeros((3, 3))
print(zeros_arr)
# Output:
# [[0. 0. 0.]
#  [0. 0. 0.]
#  [0. 0. 0.]]

ones_arr = np.ones((2, 4))
print(ones_arr)
# Output:
# [[1. 1. 1. 1.]
#  [1. 1. 1. 1.]]
```
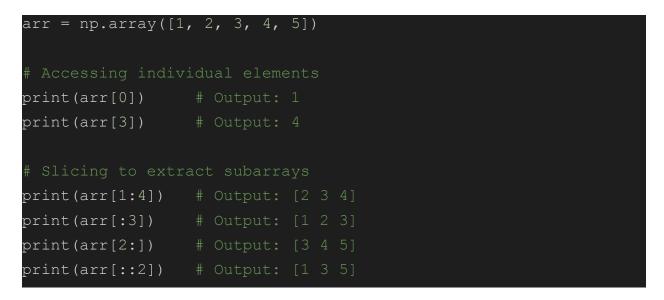
- ***Generating arrays with a range of values using `numpy.arange()`:***

You can create an array with a range of values using `np.arange()`.

```
range_arr = np.arange(1, 10, 2)
print(range_arr)   # Output: [1 3 5 7 9]
```

## 2. Indexing and Slicing NumPy Arrays:

NumPy arrays support indexing and slicing operations to access specific elements or subarrays:

```
arr = np.array([1, 2, 3, 4, 5])

# Accessing individual elements
print(arr[0])       # Output: 1
print(arr[3])       # Output: 4

# Slicing to extract subarrays
print(arr[1:4])    # Output: [2 3 4]
print(arr[:3])     # Output: [1 2 3]
print(arr[2:])     # Output: [3 4 5]
print(arr[::2])    # Output: [1 3 5]
```

In the above examples, we access individual elements by specifying the index within square brackets. We can also slice the array using a range of indices, with the syntax `[start:end:step]`. The start index is inclusive, the end index is exclusive, and the step specifies the spacing between elements.

## 3. Performing Mathematical Operations on NumPy Arrays:

NumPy arrays allow for efficient element-wise mathematical operations:

```
arr1 = np.array([1, 2, 3])
```

```
arr2 = np.array([4, 5, 6])

# Element-wise addition
addition = arr1 + arr2
print(addition)  # Output: [5 7 9]

# Element-wise multiplication
multiplication = arr1 * arr2
print(multiplication)  # Output: [4 10 18]

# Element-wise exponentiation
exponentiation = arr1 ** arr2
print(exponentiation)  # Output: [  1  32 729]

# Element-wise square root
sqrt = np.sqrt(arr1)
print(sqrt)  # Output: [1. 1.41421356 1.73205081]
```

In the above examples, the mathematical operations (`+`, `*`, `**`) are performed element-wise on the corresponding elements of the arrays. This allows for efficient computations on entire arrays without the need for explicit loops.

NumPy's array-based computations provide a powerful framework for numerical calculations and data manipulation in Python. The ability to create arrays, perform indexing and slicing operations, and execute mathematical operations efficiently makes NumPy an essential library for data analysis, scientific computing, and other numerical tasks.

# Data manipulation and analysis with pandas: reading and writing data, filtering, sorting, and aggregation

Data manipulation and analysis are core tasks in data science and analytics. The pandas library in Python provides powerful tools and functions for handling and analyzing structured data. Let's explore some key functionalities of pandas, including reading and writing data, filtering, sorting, and aggregation.

# 1. Reading and Writing Data:

pandas supports reading data from various file formats, such as CSV, Excel, SQL databases, and more. It provides functions like `read_csv()`, `read_excel()`, and `read_sql()` to import data into a pandas DataFrame. Conversely, pandas offers functions like `to_csv()`, `to_excel()`, and `to_sql()` to export data from a DataFrame.

**Example:**

```python
import pandas as pd

# Reading data from a CSV file
df = pd.read_csv("data.csv")

# Writing data to a CSV file
df.to_csv("output.csv", index=False)
```

# 2. Filtering Data:

pandas allows filtering data based on specific conditions. By applying a Boolean expression to a DataFrame or Series, we can filter rows that satisfy the given condition.

**Example:**

```python
# Filtering data based on a condition
filtered_data = df[df["Age"] > 30]
```

# 3. Sorting Data:

pandas provides the `sort_values()` function to sort data based on one or more columns. By specifying the column(s) to sort by, we can arrange the rows of the DataFrame in ascending or descending order.

**Example:**

```python
# Sorting data based on a column
sorted_data = df.sort_values("Date")
```

## 4. Aggregation:

pandas supports various aggregation operations to summarize and analyze data. The `**groupby()**` function allows grouping data based on one or more columns, followed by applying aggregation functions like `**sum()**`, `**mean()**`, `**max()**`, etc., to compute statistics for each group.

**Example:**

```python
# Grouping data and calculating the average by category
grouped_data = df.groupby("Category")["Value"].mean()
```

These are just a few examples of the wide range of data manipulation and analysis capabilities offered by pandas. With pandas, you can efficiently read and write data from different sources, filter data based on conditions, sort data, and perform aggregation operations. pandas' intuitive and expressive syntax makes it a powerful tool for handling and analyzing structured data in Python.

# Data visualization using Matplotlib: creating plots, customizing styles, and adding annotations

Data visualization is an essential aspect of data analysis, and the Matplotlib library in Python provides powerful tools for creating a wide range of visualizations. Let's explore how to create plots, customize styles, and add annotations using Matplotlib.

## 1. Creating Plots:

Matplotlib offers a variety of plot types, including line plots, scatter plots, bar plots, histograms, pie charts, and more. The `**plt.plot()**` function is commonly used to create line plots, while other specific plot types have dedicated functions like `**plt.scatter()**`, `**plt.bar()**`, etc.
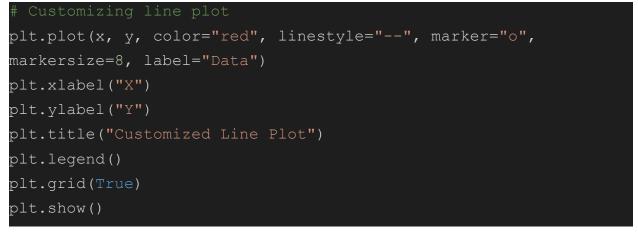
**Example:**

```python
import matplotlib.pyplot as plt

# Line plot
x = [1, 2, 3, 4, 5]
y = [10, 15, 7, 12, 9]
```

```
plt.plot(x, y)
plt.xlabel("X")
plt.ylabel("Y")
plt.title("Line Plot")
plt.show()
```

## 2. Customizing Styles:

Matplotlib provides extensive customization options to enhance the appearance of plots. You can modify elements such as colors, line styles, markers, labels, titles, and more. The `plt.plot()` function accepts additional arguments for customization, and there are many other functions like `plt.xlabel()`, `plt.ylabel()`, etc., for specific modifications.

**Example:**
```
# Customizing line plot
plt.plot(x, y, color="red", linestyle="--", marker="o",
markersize=8, label="Data")
plt.xlabel("X")
plt.ylabel("Y")
plt.title("Customized Line Plot")
plt.legend()
plt.grid(True)
plt.show()
```

## 3. Adding Annotations:

Matplotlib allows adding annotations, text, and arrows to highlight specific features or provide additional information on the plot. The `plt.text()` function is used to add text, and `plt.annotate()` is used to add annotations with arrows.

**Example:**
```
# Adding annotations
plt.plot(x, y)
plt.xlabel("X")
plt.ylabel("Y")
```

```
plt.title("Line Plot with Annotation")
plt.text(2, 13, "Important Point", fontsize=12, color="blue")
plt.annotate("Peak", xy=(3, 7), xytext=(3, 10),
arrowprops=dict(arrowstyle="->", color="red"))
plt.show()
```

These examples demonstrate how to create basic plots, customize their styles, and add annotations using Matplotlib. The library offers extensive options for customization, including colors, line styles, markers, labels, titles, and annotations, allowing you to create visually appealing and informative visualizations for data analysis. With Matplotlib's flexibility, you can tailor your plots to effectively communicate insights from your data.