∨   ** Data Toolkit (Module 8)**

## Theory Questions

### Q1. What is NumPy, and why is it widely used in Python?

Ans-> NumPy (short for Numerical Python) is a fundamental Python library for scientific computing, providing powerful tools for working with multi-dimensional arrays and matrices, and is widely used for its efficiency and versatility in data analysis, machine learning, and scientific research.

- Why is NumPy Widely Used in Python?
- Efficiency: NumPy's optimized C code and vectorized operations enable significantly faster computations, especially for large datasets and complex mathematical tasks.
- Data Analysis: NumPy is an essential tool for data analysis, providing the foundation for many data science libraries like Pandas and Scikit-learn.
- Machine Learning: NumPy's ability to efficiently handle arrays and matrices makes it a cornerstone of machine learning algorithms, which often rely on these data structures.
- Scientific Computing: NumPy is widely used in various scientific fields, including physics, biology, engineering, and finance, for simulations, modeling, and data analysis.
- Integration: NumPy integrates well with other Python libraries and allows for easy interaction with C/C++ and Fortran code.
- Vectorization: NumPy's vectorization feature allows for performing operations on entire arrays at once, eliminating the need for explicit loops and improving performance.
- Broadcasting: NumPy provides a mechanism called broadcasting, which allows operations on arrays of different shapes by automatically aligning their dimensions.

### Q2. How does broadcasting work in NumPy?

Ans-> The term broadcasting describes how NumPy treats arrays with different shapes during arithmetic operations. Subject to certain constraints, the smaller array is "broadcast" across the larger array so that they have compatible shapes.

### Q3. What is a Pandas DataFrame?

Ans-> pandas DataFrame is a way to represent and work with tabular data. It can be seen as a table that organizes data into rows and columns, making it a two-dimensional data structure. A DataFrame can be created from scratch, or you can use other data structures, like NumPy arrays.

### Q4. Explain the use of the groupby() method in Pandas?

Ans-> In Pandas, groupby() splits all the records from your data set into different categories or groups and offers you flexibility to analyze the data. The groupby() function in Pandas splits all the records from a data set into different categories or groups, offering flexibility to analyze the data by these groups.

### Q5. Why is Seaborn preferred for statistical visualizations?

Ans-> One of the main advantages of Seaborn over Matplotlib is its ease of use. Seaborn is built on top of Matplotlib and provides a higher-level interface for creating statistical graphics. This means that Seaborn requires less code to create complex visualizations compared to Matplotlib.

### Q6. What are difference between NumPy arrays and Python lists?

Ans-> The biggest difference is NumPy arrays use fewer resources than Python lists, which becomes important when storing a large amount of data. If you're working with thousands of elements, Python lists will be fine for most purposes.

- Python Lists: Python lists are fundamental data structures in Python, offering flexibility and versatility. Here are some key points:

- Dynamic Typing: Lists can accommodate elements of different data types within the same list.

- Size Flexibility: Lists can dynamically grow or shrink in size as elements are added or removed.

- Versatility: Lists support a variety of built-in operations and methods for manipulation.

- Slower Performance: Lists tend to be slower compared to NumPy arrays, especially for large datasets, due to the overhead of dynamic typing and flexibility. Less Mathematical Functionality: Lists lack built-in support for mathematical operations on entire arrays of data.

- NumPy Arrays: NumPy arrays are optimized data structures for numerical computation and data manipulation. Here's what distinguishes them:

- Homogeneous Data Types: NumPy arrays require all elements to be of the same data type, enabling efficient storage and operations.

- Fixed Size: Arrays have a fixed size upon creation, enhancing memory usage and computational efficiency.
- Optimized Operations: NumPy provides a plethora of optimized mathematical operations and functions that can be applied directly to arrays, making them * ideal for numerical tasks. Better Performance: NumPy arrays offer significantly better performance compared to Python lists, especially for large datasets, due to their fixed size and homogeneous data type.
- Memory Efficiency: NumPy arrays are more memory efficient than Python lists, particularly for large datasets, because they store data in contiguous memory blocks.
- Broadcasting: NumPy arrays support broadcasting, enabling operations between arrays of different shapes and sizes without explicit iteration.

## Q7. What is a heatmap, and when should it be used?

Ans-> A heatmap is a data visualization technique that uses color variations to represent the magnitude of values in a dataset. It's essentially a 2D grid where each cell's color corresponds to the value at that location. Heatmaps are useful for quickly identifying patterns, trends, and areas of high and low concentration within a dataset.

  - When to use heatmaps:

- Website analytics: To understand user behavior, such as where users click, scroll, or hover on a page.
- Data analysis: To visualize complex data, identify correlations between variables, and see distributions of values.
- Geographic data: To map data like population density, temperature, or crime rates.
- Business performance: To track performance of different business units or departments.
- Correlation analysis: To visualize relationships between multiple variables.
- Earth science: To visualize data like weather patterns, climate change, or satellite imagery.

## Q8. What does the term "vectorized operation" mean in NumPy?

Ans-> Vectorized operations, on the other hand, take advantage of low-level optimizations implemented in libraries like NumPy (which pandas is built on). These operations apply a function to an entire array (or DataFrame/Series) at once, leveraging highly efficient C and Fortran code under the hood.

## Q9. How does matplotlib differ from Plotly?

Ans-> Both Matplotlib and Plotly have their distinct advantages. Matplotlib offers unmatched control and is perfect for creating static, publication-quality visuals. Plotly, on the other hand, excels in interactivity and ease of use, making it ideal for dynamic presentations and web applications.

## Q10. What is significance of hierarchical indexing in Pandas?

Ans-> Hierarchical Indexing, also known as MultiIndexing, is a powerful feature in Pandas that allows you to have multiple levels of indexing on an axis (row or column). This capability is particularly useful when dealing with high-dimensional data.

## Q11. What is the role of Seaborn's pairplot() in function?

Ans-> Seaborn's pairplot() function creates a grid of plots that visually display pairwise relationships between variables in a dataset. It effectively summarizes data distributions and correlations in a single, comprehensive figure. The diagonal of the plot shows histograms or kernel density estimates, while the off-diagonal plots display scatterplots representing the relationships between pairs of variables.

## Q12. What is the purpose of the describe() function in Pandas?

Ans-> The purpose of the describe() function in Pandas is to generate descriptive statistics for a DataFrame or Series. It provides a summary of the data, including the count, mean, standard deviation, minimum, 25th percentile (Q1), median (50th percentile or Q2), 75th percentile (Q3), and maximum values for each numerical column. For non-numerical columns, it also provides statistics like unique values, top value, and frequency.

## Q13. Why is handling missing data important in Pandas?

Ans-> Handling missing data in Pandas is crucial because it ensures data integrity, avoids skewed analysis, and improves the reliability of machine learning models. Missing values, often represented as NaN (Not a Number) in Pandas, can lead to incorrect calculations, misleading conclusions, and reduced statistical power.

- Here's why it's important:
- Data Integrity: Missing values can distort calculations and lead to inaccurate representations of the data.
- Bias and Accuracy: Without proper handling, missing data can introduce biases and inaccuracies into the analysis, impacting everything from data visualization to machine learning models.
- Reliability of Models: Machine learning models can struggle with datasets containing missing data, leading to unreliable predictions.

- Effective Analysis: Missing data can skew analysis results, making it difficult to draw accurate conclusions.
- Statistical Power: Missing values can reduce the power and validity of statistical models.

## Q14. What are the benefits of using plotly for data visualization?

Ans-> Flexibility and Customization You can also customize practically every element of the chart, from the color scheme to the hover tooltips. This means that if you're working on a project where aesthetics are important, Plotly gives you full control over the look and feel of your visualizations.

## Q15. How does NumPy handle multimensional arrays?

Ans-> For numpy, all data in any multi-dimensional array is actually stored in memory as a long 1D array (we will get back to this in the master lecture). The number of dimensions and shape of an array is actually only used to structure the data access in a certain way.

## Q16. What is the role of Bokeh in data visulization?

Ans-> Bokeh is a Python library primarily used for creating interactive and web-based data visualizations. It allows users to build beautiful and dynamic charts, plots, and dashboards that can be displayed in web browsers, including within applications and Jupyter notebooks. Bokeh distinguishes itself from other Python visualization libraries like Matplotlib or Seaborn by focusing on interactive elements and browser output.

## Q17. Explain the difference between apply() and map() in Pandas?

Ans-> 1 Map Function The map method in pandas is a Series method used primarily for substituting each value in a Series with another value. Imagine you have a Series of categories that you want to convert to numerical codes. You can use map with a dictionary where keys are current items and values are the intended substitutions. This method is particularly useful for element-wise transformations and is limited to working with Series objects only, not DataFrames.

2 Apply Method Moving on to the apply method, this is more versatile as it works on both Series and DataFrame objects. When used on a Series, it behaves similarly to map , applying a function to each element. However, when apply is used on a DataFrame, it can operate across rows or columns. For example, if you want to calculate the mean of each column, you can use apply with the axis parameter set to 0. This flexibility makes apply a go-to for many row or column-wise operations.

## Q18. What are some advanced features of NumPy?

Ans-> NumPy's advanced features include structured arrays, structured record arrays, broadcasting, fancy indexing, and universal functions (ufuncs), which provide powerful and efficient ways to work with numerical data.

- Here's a more detailed look:

1. Structured and Record Arrays: Structured Arrays: These arrays can hold different data types in different columns, allowing for more flexible data organization. Record Arrays: These are a specific type of structured array that provides a way to access elements using field names, making data access more intuitive.
2. Broadcasting: Broadcasting allows NumPy to perform operations on arrays of different shapes. It automatically expands the smaller array's dimensions to match the larger array, making operations simpler and more efficient. This feature is particularly useful when working with arrays of uneven shapes.
3. Fancy Indexing: Fancy indexing provides a way to access array elements using a list or array of indices. This is more powerful than simple slicing and allows for more complex indexing patterns.
4. Universal Functions (ufuncs): Ufuncs are functions that operate on NumPy arrays element-wise. They are highly optimized and can perform mathematical, logical, and other operations on arrays efficiently. Examples of ufuncs include add, subtract, multiply, divide, and power.
5. Other Advanced Features: Aggregation: NumPy provides functions for aggregating data, such as calculating sums, means, and standard deviations. Masking: Masking allows you to selectively exclude elements from an array based on certain conditions, enabling more targeted operations. Sorting: NumPy provides functions for sorting arrays in different ways, including sorting by specific columns or axes. Linear Algebra: NumPy includes a powerful

## Q19. How does Pandas simplify time series analysis?

Ans-> Combining the ease of use of dateutil and datetime modules and the vectorized interface and efficient storage of NumPy's datetime64, pandas provides a Timestamp object. The library then makes a DatetimeIndex from these Timestamp objects to index a DataFrame or Series.

## Q20. What is the role of a pivot table in Pandas?

Ans-> The pivot() function is an incredibly useful tool for transforming and summarizing data. It allows you to restructure a DataFrame by turning rows into columns and columns into rows based on a specified index column, a specified columns column, and a specified values column.

Pandas pivot tables work in a very similar way to those found in spreadsheet tools such as Microsoft Excel. The pivot table function takes in a data frame and the parameters detailing the shape you want the data to take. Then it outputs summarized data in the form of a pivot table.

## Q21. Why is NumPy's array slicing faster than Python's list slicing?

Ans-> Efficient Memory Usage: NumPy arrays are more memory-efficient than lists because they store elements of the same data type in contiguous memory locations. Lists, by contrast, can hold objects of different types, requiring more overhead.

## Q22. What are some common use cases of Seaborn?

Ans-> Seaborn is a library for creating statistical graphics in Python. It is based on Matplotlib and integrates with Pandas data structures. This library is as powerful as Matplotlib but brings simplicity and unique features. It allows for quick data exploration and understanding.

Scatter Plot. A scatter plot is used to visualize the relationship between two variables

- Line Plot. A line plot is used to visualize the trend of a variable over * * time
- Histogram
- Box Plot
- Violin Plot
- Heatmap
- Pairplot

## Practical Questions

## Q1. How do you create a 2D Numpy array and calculate the sum of each row?

```
# prompt: How do you create a 2D Numpy array and calculate the sum of each row

import numpy as np

# Create a 2D NumPy array
arr_2d = np.array([[1, 2, 3],
                   [4, 5, 6],
                   [7, 8, 9]])

# Calculate the sum of each row
row_sums = np.sum(arr_2d, axis=1)

# Print the array and row sums
print("2D Array:\n", arr_2d)
print("\nSum of each row:\n", row_sums)
```

```
2D Array:
 [[1 2 3]
 [4 5 6]
 [7 8 9]]

Sum of each row:
 [ 6 15 24]
```

## Q2. Write a Pandas script to find the mean of a specific column in a DataFrame?

```
# prompt: Write a Pandas script to find the mean of a specific column in a DataFrame

import pandas as pd

# Sample DataFrame (replace with your actual data)
data = {'col1': [1, 2, 3, 4, 5],
        'col2': [6, 7, 8, 9, 10],
        'col3': [11, 12, 13, 14, 15]}
df = pd.DataFrame(data)

# Specify the column name
column_name = 'col2'  # Replace with the actual column name

# Calculate the mean of the specified column
mean_value = df[column_name].mean()

# Print the mean
print(f"The mean of column '{column_name}' is: {mean_value}")
```

```
The mean of column 'col2' is: 8.0
```

## ⌄ Q3. Create a scatter plot using Matplotlib?

```
# prompt: Create a scatter plot using Matplotlib

import matplotlib.pyplot as plt

# Sample data (replace with your actual data)
x = [1, 2, 3, 4, 5]
y = [2, 4, 1, 5, 3]

# Create the scatter plot
plt.scatter(x, y)

# Add labels and title
plt.xlabel("X-axis")
plt.ylabel("Y-axis")
plt.title("Scatter Plot")

# Display the plot
plt.show()
```
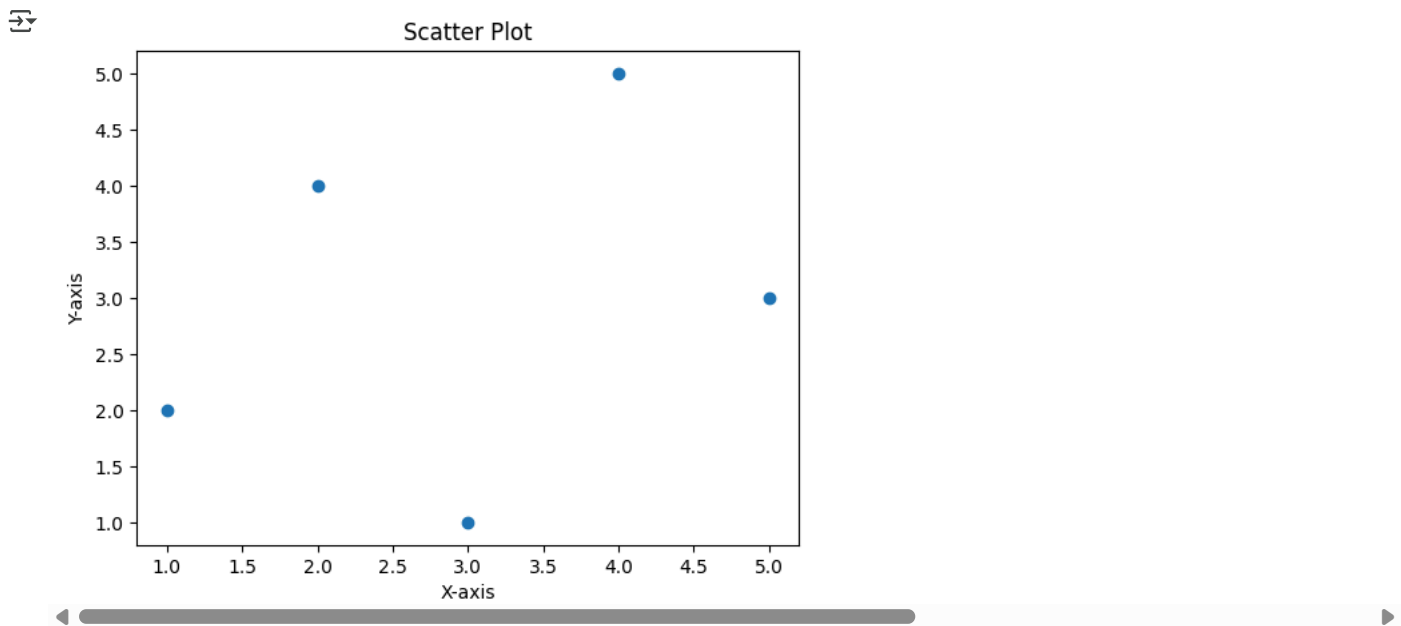


## ⌄ Q4. How do you calculate the correlation matrix using Seaborn and visulize it with a heatmap?

```
# prompt: How do you calculate the correlation matrix using Seaborn and visulize it with a heatmap

import seaborn as sns
import matplotlib.pyplot as plt

# Sample DataFrame (replace with your actual data)
data = {'col1': [1, 2, 3, 4, 5],
        'col2': [6, 7, 8, 9, 10],
        'col3': [11, 12, 13, 14, 15]}
df = pd.DataFrame(data)

# Calculate the correlation matrix
correlation_matrix = df.corr()

# Create the heatmap
plt.figure(figsize=(8, 6))  # Adjust figure size as needed
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
plt.title('Correlation Matrix Heatmap')
plt.show()
```
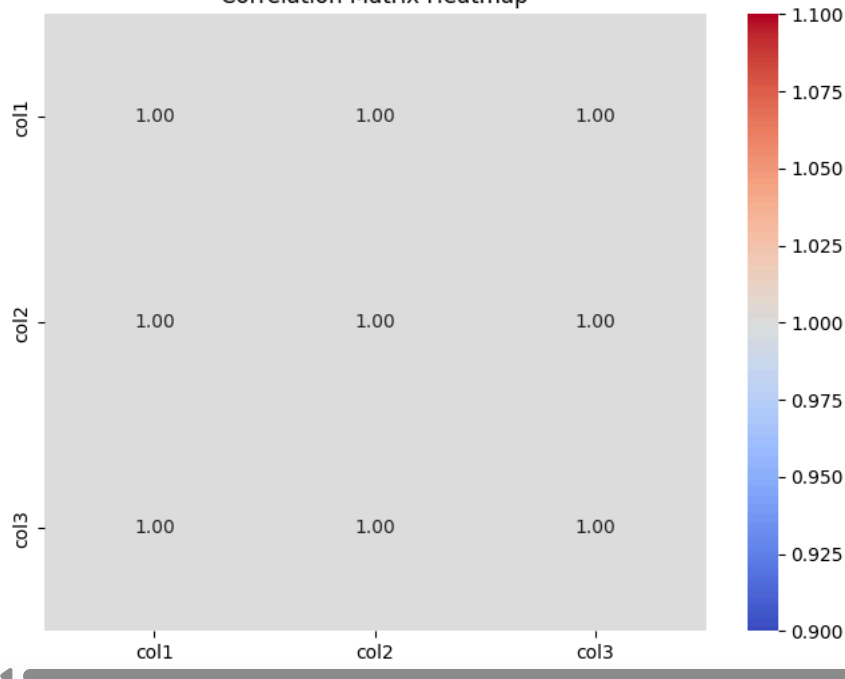
## Q5. Generate a bar plot using Matplotlib

```
# prompt: generate a bar plot using matplotlib

# Sample data (replace with your actual data)
categories = ['A', 'B', 'C', 'D', 'E']
values = [10, 15, 7, 12, 9]

# Create the bar plot
plt.bar(categories, values)

# Customize the plot (optional)
plt.xlabel("Categories")
plt.ylabel("Values")
plt.title("Bar Plot")

# Display the plot
plt.show()
```
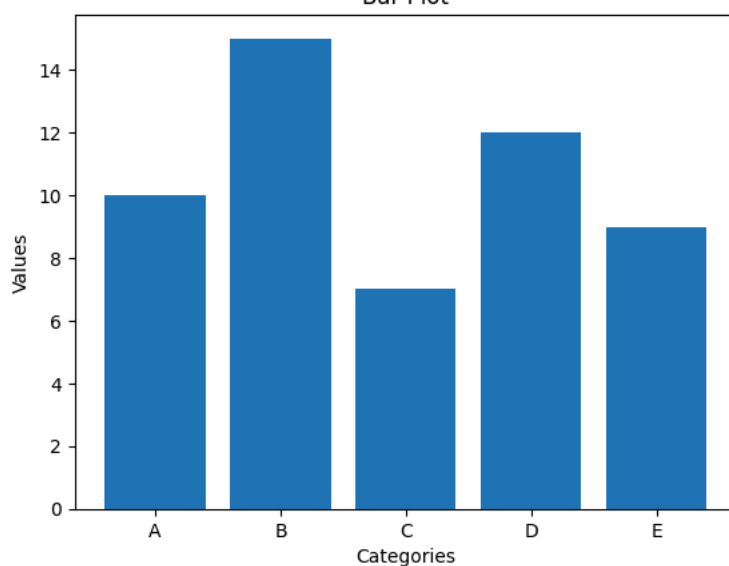


## Q6. Create a DataFrame and add a new column based on an existing column

```
# prompt: Create a DataFrame and add a new column based on an existing column
```

```python
import pandas as pd

# Create a sample DataFrame
data = {'col1': [1, 2, 3, 4, 5]}
df = pd.DataFrame(data)

# Add a new column 'col2' based on 'col1'
df['col2'] = df['col1'] * 2

# Print the DataFrame
df
```

|   | col1 | col2 |
|---|------|------|
| 0 | 1    | 2    |
| 1 | 2    | 4    |
| 2 | 3    | 6    |
| 3 | 4    | 8    |
| 4 | 5    | 10   |

Next steps:  [ Generate code with df ]  [ ⊙ View recommended plots ]  [ New interactive sheet ]
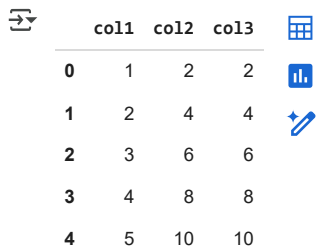
## ⌄ Q7. Write a program to perform element-wise multiplication of Two NumPy arrays

```python
# prompt: Create a DataFrame and add a new column based on an existing column

# Assuming you have a DataFrame named 'df' with a column named 'col1'
# (as shown in the previous code example)

# Add a new column 'col3' where each element is twice the value in 'col1'
df['col3'] = df['col1'] * 2

# Print the updated DataFrame
df
```

|   | col1 | col2 | col3 |
|---|------|------|------|
| 0 | 1    | 2    | 2    |
| 1 | 2    | 4    | 4    |
| 2 | 3    | 6    | 6    |
| 3 | 4    | 8    | 8    |
| 4 | 5    | 10   | 10   |

Next steps:  [ Generate code with df ]  [ ⊙ View recommended plots ]  [ New interactive sheet ]

## ⌄ Q8. Create a line plot with multiple lines using Matplotlib

```python
# prompt: Create a line plot with multiple lines using Matplotlib

# Sample data (replace with your actual data)
x = np.linspace(0, 10, 100)  # 100 evenly spaced points from 0 to 10
y1 = np.sin(x)
y2 = np.cos(x)
y3 = x**2

# Create the plot
plt.plot(x, y1, label='sin(x)')
plt.plot(x, y2, label='cos(x)')
plt.plot(x, y3, label='x^2')

# Add labels and title
plt.xlabel("X-axis")
plt.ylabel("Y-axis")
plt.title("Multiple Line Plot")

# Add a legend
plt.legend()

# Display the plot
```
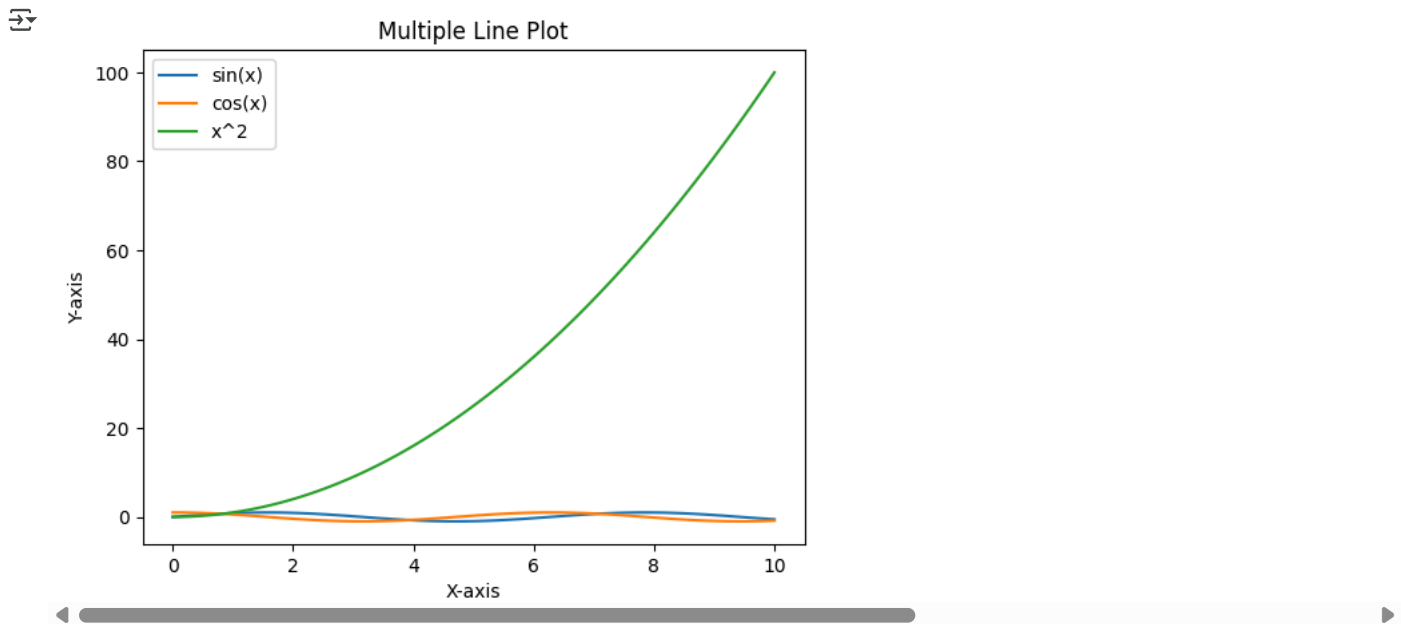
```
plt.show()
```


Multiple Line Plot

## Q9. Generate a Pandas DataFrame and filter rows where a column value is greater than a threshold

```
# prompt: Generate a Pandas DataFrame and filter rows where a column value is greater than a threshold

import pandas as pd

# Create a sample DataFrame
data = {'col1': [1, 2, 3, 4, 5],
        'col2': [6, 7, 8, 9, 10]}
df = pd.DataFrame(data)

# Set the threshold value
threshold = 3

# Filter rows where 'col1' is greater than the threshold
filtered_df = df[df['col1'] > threshold]

# Print the filtered DataFrame
filtered_df
```

|   | col1 | col2 |
|---|------|------|
| 3 | 4    | 9    |
| 4 | 5    | 10   |

Next steps: ( Generate code with `filtered_df` ) ( 👁 View recommended plots ) ( New interactive sheet )

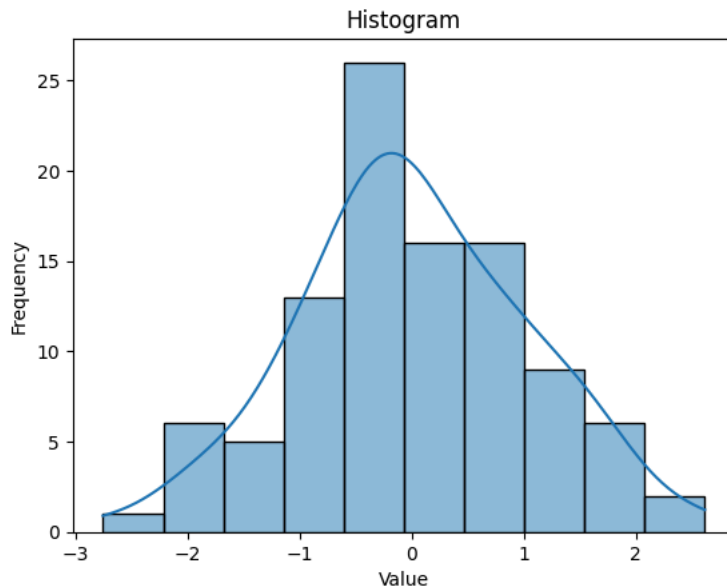## Q10. Create a histogram using Seaborn to visualize a distribution

```
# prompt: Create a histogram using Seaborn to visualize a distribution

# Sample data (replace with your actual data)
data = np.random.randn(100)  # 100 random numbers from a standard normal distribution

# Create the histogram
sns.histplot(data, kde=True)  # kde=True adds a kernel density estimate curve

# Customize the plot (optional)
plt.xlabel("Value")
plt.ylabel("Frequency")
plt.title("Histogram")

# Display the plot
plt.show()
```

## Q11. Perform Matrix multiplication using NumPy

```python
# prompt: Perform Matrix multiplication using NumPy

# Create two matrices
matrix1 = np.array([[1, 2], [3, 4]])
matrix2 = np.array([[5, 6], [7, 8]])

# Perform matrix multiplication
result_matrix = np.dot(matrix1, matrix2)

# Print the result
print("Matrix 1:\n", matrix1)
print("\nMatrix 2:\n", matrix2)
print("\nResultant Matrix:\n", result_matrix)
```

```
Matrix 1:
 [[1 2]
 [3 4]]

Matrix 2:
 [[5 6]
 [7 8]]

Resultant Matrix:
 [[19 22]
 [43 50]]
```

## Q12. Use Pandas to load a CSV file and display its first 5 rows

```python
# prompt: Use Pandas to load a CSV file and display its first 5 rows

import pandas as pd

# Assuming the CSV file is named 'your_file.csv' and is in the current directory
# Replace 'your_file.csv' with the actual file name and path if needed
try:
  df = pd.read_csv('your_file.csv')
  print(df.head())
except FileNotFoundError:
  print("Error: 'your_file.csv' not found. Please check the file name and path.")
```

```
Error: 'your_file.csv' not found. Please check the file name and path.
```

## Q13. Create a 3D Scatter plot using Plotly

```python
# prompt: create a 3d scatter plot using plotly
```

```
import plotly.graph_objects as go
import numpy as np

# Sample data (replace with your actual data)
x = np.random.rand(100)
y = np.random.rand(100)
z = np.random.rand(100)

# Create the 3D scatter plot
fig = go.Figure(data=[go.Scatter3d(x=x, y=y, z=z, mode='markers')])

# Customize the plot (optional)
fig.update_layout(scene=dict(
    xaxis_title='X-axis',
    yaxis_title='Y-axis',
    zaxis_title='Z-axis'
))

# Display the plot
fig.show()
```