

AWORTWE ENOCK

[enockawor@gmail.com](mailto:enockawor@gmail.com)

+233 245227067

## **Matplotlib and Pandas**

### **Library Overview**

**Matplotlib** is a popular Python package for producing static, animated, and interactive visualizations. It offers a thorough foundation for creating many different types of plots and charts, such as bar charts, histograms, scatter plots, line graphs, and more. Matplotlib's enormous customization possibilities are one of its most notable features; they let users adjust nearly every component of their plots to get the exact look and feel they want. It is therefore especially appropriate for generating figures of publishing quality. Matplotlib is a versatile and customizable plot tool that allows users to create various plot types, control line styles, colors, and fonts. It seamlessly integrates with NumPy and pandas, supports complex layouts with multiple subplots and figures, and can produce interactive plots using backends like Tkinter and wxPython. Typical use cases include scientific research, data analysis, engineering, and finance, where charts help visualize data distributions, trends, and track stock prices and financial metrics.

**Pandas** is an open-source data analysis and manipulation package for Python that is strong and adaptable. It offers the operations and data structures required to easily handle structured data. Built on top of NumPy arrays for effective numerical operations, its main data structures, Series (1-dimensional) and DataFrame (2-dimensional), are intended to handle a broad range of data types. Pandas' ease of handling and processing of huge datasets makes it a popular tool in data science, machine learning, and data analysis. The software offers data structures, data alignment, integration with Python libraries like NumPy, Matplotlib, and Scikit-learn, data manipulation, and time series tools for handling labeled data, reshaping, merging, and aggregating data. Its typical use cases are data cleaning which involves handling missing data, transforming data types, and preparing data for analysis. EDA (Exploratory Data Analysis) which summarizes statistics, data transformation applies operations across datasets, time series analysis visualizes trends, and data aggregation groups and summarizes data for insights.

### **Graphic Types**

#### **1. Matplotlib**

A large range of graph types are available for data visualization using Matplotlib. This is a list of some of the more popular graph types that are supported by Matplotlib, complete with a synopsis, possible applications, and basic code samples:

##### **a. Line Plot**

A line plot is a visual representation of data, consisting of markers connected by straight lines, used to display trends over time or continuous data.

Code Example:

```
import matplotlib.pyplot as plt

x = [1, 2, 3, 4, 5]
y = [2, 3, 5, 7, 11]

plt.plot(x, y)

plt.title('Line Plot')

plt.xlabel('X-axis')

plt.ylabel('Y-axis')

plt.show()
```

b. Scatter Plot

A scatter plot uses dots to represent values for two numeric variables, indicating individual data points on the horizontal and vertical axis, aiding in identifying relationships or correlations.

Code Example:

```
import matplotlib.pyplot as plt

x = [1, 2, 3, 4, 5]
y = [2, 3, 5, 7, 11]

plt.scatter(x, y)

plt.title('Scatter Plot')

plt.xlabel('X-axis')

plt.ylabel('Y-axis')

plt.show()
```

c. Bar Chart

A bar chart is a visual tool used to display categorical data, with rectangular bars representing values proportional to their heights or lengths.

```
import matplotlib.pyplot as plt

categories = ['A', 'B', 'C', 'D', 'E']
values = [5, 7, 3, 8, 4]

plt.bar(categories, values)

plt.title('Bar Chart')

plt.xlabel('Categories')

plt.ylabel('Values')

plt.show()
```

d. Histogram

A histogram is a visual representation of numerical data distribution, typically used for frequency distributions, and is particularly useful for identifying outliers in a dataset.

```
import matplotlib.pyplot as plt

data = [1, 2, 2, 3, 3, 3, 4, 4, 4, 4, 5, 5, 6]

plt.hist(data, bins=5)

plt.title('Histogram')

plt.xlabel('Value')

plt.ylabel('Frequency')

plt.show()
```

e. Pie Chart

A pie chart is a circular statistical graphic divided into slices, used to illustrate numerical proportions and relative proportions of a whole.

```
import matplotlib.pyplot as plt

labels = ['A', 'B', 'C', 'D']
sizes = [15, 30, 45, 10]

plt.pie(sizes, labels=labels,
        autopct='%1.1f%%')

plt.title('Pie Chart')

plt.show()
```

## 2. Pandas

Using Matplotlib as its backend, Pandas offers a variety of simple techniques for data visualization. This is a list of some of the more popular graph types that Pandas supports, complete with a synopsis, possible applications, and basic code samples.

### a. Line Plot

A line plot is a visual representation of data, consisting of markers connected by straight lines, used to display trends over time or continuous data.

```
import pandas as pd

data = {'x': [1, 2, 3, 4, 5], 'y': [2, 3, 5, 7, 11]}

df = pd.DataFrame(data)

df.plot.line(x='x', y='y', title='Line Plot')

plt.xlabel('X-axis')

plt.ylabel('Y-axis')

plt.show()
```

### b. Scatter Plot

A scatter plot uses dots to represent values for two numeric variables, indicating individual data points on the horizontal and vertical axis, aiding in identifying relationships or correlations.

```
import pandas as pd

data = {'x': [1, 2, 3, 4, 5], 'y': [2, 3, 5, 7, 11]}

df = pd.DataFrame(data)

df.plot.scatter(x='x', y='y',
               title='Scatter Plot')

plt.xlabel('X-axis')

plt.ylabel('Y-axis')

plt.show()
```

c. Bar Chart

A bar chart is a visual tool used to display categorical data, with rectangular bars representing values proportional to their heights or lengths.

```
import pandas as pd

data = {'Category': ['A', 'B', 'C', 'D', 'E'], 'Value': [5, 7, 3, 8, 4]}

df = pd.DataFrame(data)

df.plot.bar(x='Category', y='Value', title='Bar Chart')

plt.xlabel('Categories')

plt.ylabel('Values')

plt.show()
```

d. Histogram

A histogram is a visual representation of numerical data distribution, typically used for frequency distributions, and is particularly useful for identifying outliers in a dataset.

```
import pandas as pd

data = {'Value': [1, 2, 2, 3, 3, 3, 4, 4, 4, 4, 5, 5, 6]}

df = pd.DataFrame(data)

df.plot.hist(y='Value', bins=5,
title='Histogram')

plt.xlabel('Value')

plt.ylabel('Frequency')

plt.show()
```

e. Pie Chart

A pie chart is a circular statistical graphic divided into slices, used to illustrate numerical proportions and relative proportions of a whole.

```
import pandas as pd

data = {'Category': ['A', 'B', 'C', 'D'],
'Value': [15, 30, 45, 10]}

df = pd.DataFrame(data)

df.set_index('Category').plot.pie(y='Value', autopct='%1.1f%%', title='Pie Chart')

plt.ylabel('')

plt.show()
```

## **Comparisons**

While both Matplotlib and Pandas are effective Python tools for data visualization, they have different advantages and disadvantages. Here is a thorough comparison that takes into account factors like performance with big datasets, customization possibilities, and ease of use:

1. Ease of Use

Matplotlib offers a wide range of plotting functions and fine-grained control over plot elements, with extensive documentation and a large user community. It is suitable for creating complex and publication-quality plots. However, it has a steeper learning curve and requires more code for simple plots. Pandas simplifies the plotting process by integrating with DataFrames, but is limited to straightforward plots and less suitable for complex visualizations.

2. Customization Options

Matplotlib is highly customizable, offering control over various plot aspects like colors, labels, scales, and styles. It supports advanced plotting features like subplots and annotations. Pandas inherits customization options from Matplotlib, making basic customization straightforward. It's suitable for quick and moderately customized visualizations but has limited advanced customization and less control over fine details.

3. Interactivity

Matplotlib supports interactive plots through backends like Tkinter and Qt and integrates with Jupyter Notebooks. It has advanced interactivity features like Plotly or Bokeh. Pandas can generate interactive plots within Jupyter Notebooks and integrates with Plotly for enhanced interactivity.

4. Performance with large dataset

Matplotlib is efficient for moderate-sized datasets and can handle large datasets with proper optimization. Pandas is efficient for quick data plotting from DataFrames, suitable for moderately large datasets. However, performance may degrade with large datasets and slow plotting may occur.

## References

- Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*, 9(3), 90-95. doi:10.1109/MCSE.2007.55
- Droettboom, M., Caswell, T. A., Hunter, J., Elson, P., Firing, E., Nielsen, J. H., ... & West, C. K. (2018). Matplotlib: Version 3.0.2. Zenodo. doi:10.5281/zenodo.1482090
- McKinney, W. (2010). Data Structures for Statistical Computing in Python. *Proceedings of the 9th Python in Science Conference*, 51-56. doi:10.25080/Majora-92bf1922-00a
- Reback, J., McKinney, W., jbrockmendel, Van den Bossche, J., Augspurger, T., Cloud, P., ... & Sinhrks. (2020). pandas-dev/pandas: Pandas 1.0.3. Zenodo. doi:10.5281/zenodo.3715232