

# Data Analytics and Visualization with Python

## Q. Explain about different libraries for Data Analytics in Python.

Python offers a **rich ecosystem of libraries** for data analytics, covering everything from **data manipulation and visualization to machine learning and big data processing**. Below is a detailed breakdown of the most important libraries:

---

### 1. Data Manipulation & Processing Libraries

**1.1 Pandas (Python Data Analysis Library):** Pandas is a powerful data manipulation and analysis library built on top of NumPy. It provides DataFrame and Series objects for handling structured data, similar to SQL tables or Excel spreadsheets.

#### Key Features:

- Handles tabular data (CSV, Excel, SQL) easily.
- Provides functions for data cleaning, filtering, merging, and grouping.
- Supports time series analysis.

#### Example Usage:

```
import pandas as pd
df = pd.read_csv("data.csv") # Read a CSV file
df.info() # Get dataset summary
df.describe() # Get basic statistics
df.dropna(inplace=True) # Remove missing values
df["new_column"] = df["old_column"] * 2 # Create a new column
```

**Best for:** Data preprocessing, analysis, and handling large datasets.

---

**1.2 NumPy (Numerical Python):** NumPy is the foundation of numerical computing in Python. It provides fast, efficient operations on multi-dimensional arrays and matrices.

#### Key Features:

- Optimized for vectorized operations (faster than Python lists).
- Provides functions for mathematical operations (mean, median, standard deviation, etc.).
- Works as the base for Pandas, SciPy, Scikit-Learn, and TensorFlow.

#### Example Usage:

```
import numpy as np
```

```
arr = np.array([1, 2, 3, 4, 5]) # Create an array
print(arr.mean()) # Calculate mean
print(arr[1:4]) # Slicing
```

**Best for:** Fast mathematical computations and working with numerical data.

---

## 2. Data Visualization Libraries

**2.1 Matplotlib:** Matplotlib is the most widely used library for static plots in Python. It allows users to create line plots, bar charts, histograms, scatter plots, and more.

### Key Features:

- Highly customizable but requires more code than Seaborn.
- Supports 2D & 3D plotting.
- Works well with Pandas and NumPy.

### Example Usage:

```
import matplotlib.pyplot as plt
plt.plot([1, 2, 3, 4], [10, 20, 30, 40]) # Line plot
plt.xlabel("X-axis")
plt.ylabel("Y-axis")
plt.title("Simple Line Graph")
plt.show()
```

**Best for:** Detailed control over plots but requires more manual customization.

---

**2.2 Seaborn:** Seaborn is built on Matplotlib and is designed for statistical data visualization. It simplifies the creation of beautiful, informative graphs.

### Key Features:

- Easier to use than Matplotlib with built-in themes.
- Supports categorical plots, heatmaps, violin plots, and regression plots.
- Works well with Pandas DataFrames.

### Example Usage:

```
import seaborn as sns
sns.histplot(df["column_name"], kde=True) # Histogram with density curve
sns.boxplot(x=df["category"], y=df["values"]) # Boxplot
```

**Best for:** Quick and stylish statistical plots.

---

**2.3 Plotly:** Plotly is an interactive visualization library that allows zooming, hovering, and dynamic updates, making it perfect for dashboards and web applications.

**Key Features:**

- Creates interactive graphs for web applications.
- Supports 3D plots and real-time streaming data.
- Works with Dash framework for building dashboards.

**Example Usage:**

```
import plotly.express as px  
  
fig = px.scatter(df, x="column1", y="column2", color="category")  
  
fig.show()
```

**Best for:** Interactive visualizations and dashboards.

---

### 3. Machine Learning & Statistical Analysis Libraries

**3.1 Scikit-Learn:** Scikit-Learn is the most popular machine learning library for implementing algorithms like linear regression, decision trees, and clustering.

**Key Features:**

- Supports classification, regression, clustering, and dimensionality reduction.
- Provides model selection tools (train-test split, cross-validation, hyperparameter tuning).

**Example Usage:**

```
from sklearn.linear_model import LinearRegression  
  
model = LinearRegression()  
  
model.fit(X_train, y_train)  
  
y_pred = model.predict(X_test)
```

**Best for:** Traditional ML models (classification, regression, clustering).

---

**3.2 Statsmodels:** Statsmodels is used for advanced statistical modeling, such as linear regression, hypothesis testing, and time series analysis.

**Key Features:**

- Includes t-tests, ANOVA, and Generalized Linear Models (GLM).
- More detailed statistical output than Scikit-Learn.

**Example Usage:**

```
import statsmodels.api as sm
```

```
model = sm.OLS(y, X).fit()
print(model.summary()) # Detailed statistical summary
```

**Best for:** Hypothesis testing, statistical modeling.

## Q. Explain process of creating Histogram using Python

A **histogram** is a graphical representation of the **distribution of numerical data**. It groups data into **bins (intervals)** and represents the frequency of data points within each bin.

In Python, we can create histograms using the following libraries:

1. **Matplotlib** – Basic plotting library
2. **Seaborn** – Advanced statistical visualization
3. **Pandas** – Quick histogram from DataFrame

---

### 1. Creating a Histogram Using Matplotlib

**Matplotlib's hist() function** is the most commonly used method to create a histogram.

**Steps to Create a Histogram using Matplotlib:**

1. Import **Matplotlib** and **NumPy**.
2. Generate or load a dataset.
3. Use `plt.hist()` to plot the histogram.
4. Customize the number of **bins**, colors, labels, and titles.
5. Display the plot using `plt.show()`.

**Example: Basic Histogram**

```
import numpy as np
import matplotlib.pyplot as plt

# Generate random data
data = np.random.randn(1000) # 1000 random values

# Create histogram
plt.hist(data, bins=30, color="skyblue", edgecolor="black")
```

```
# Add labels and title
plt.xlabel("Value")
plt.ylabel("Frequency")
plt.title("Histogram of Random Data")
```

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

## 2. Creating a Histogram Using Seaborn

**Seaborn's histplot() function** makes it easier to create attractive statistical plots.

### Steps to Create a Histogram using Seaborn:

1. Import **Seaborn** and **Matplotlib**.
2. Load the dataset.
3. Use `sns.histplot()` to create the histogram.
4. Customize the number of bins, colors, and other properties.

### Example: Histogram with KDE (Kernel Density Estimate)

```
import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt

# Generate random data
data = np.random.randn(1000)

# Create histogram using Seaborn
sns.histplot(data, bins=30, kde=True, color="green")

# Add labels and title
plt.xlabel("Value")
plt.ylabel("Frequency")
plt.title("Seaborn Histogram with KDE")

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