

VISUALIZATION

Visualization is a way of representing. Data visualization is representing of data in a graphical or pictorial format for better understanding of data. Visualization gives a good idea of data and the trends in it.

In python, data visualization has multiple libraries like matplotlib, seaborn, plotly etc.

Let us dive deeper into matplotlib and seaborn.

MATPLOTLIB

Matplotlib, introduced in 2002 by John Hunter. It is for plotting 2D plots which is built on numpy. A plot of matplotlib contains:

1 Figure

2 Axes

3 Axis

4 Artists

Figure

Figure is the container in which plots are present. It can have a single plot or multiple plots.

Axes

Axes are the plots, and they are artist attached to the figure. A plot can have 2 or 3 axes. `set_xlabel()`, `set_ylabel()` are the functions used to set the labels of x and y respectively.

Axis

It is responsible for generating ticks or the limits on the axes.

Artists

The whatever content visible in a figure are the artists.

PYPLLOT

Pyplot is a sub library of matplotlib where all the utilities lie under. It has different types of plots including bar graphs, scatter plots, pie charts, histograms, area charts.

We import pyplot from matplotlib as following:

```
1 import matplotlib.pyplot as plt
```

We also import numpy as a support for the matplotlib :

```
import numpy as np
```

Some of the plots in matplotlib:

LINE CHART

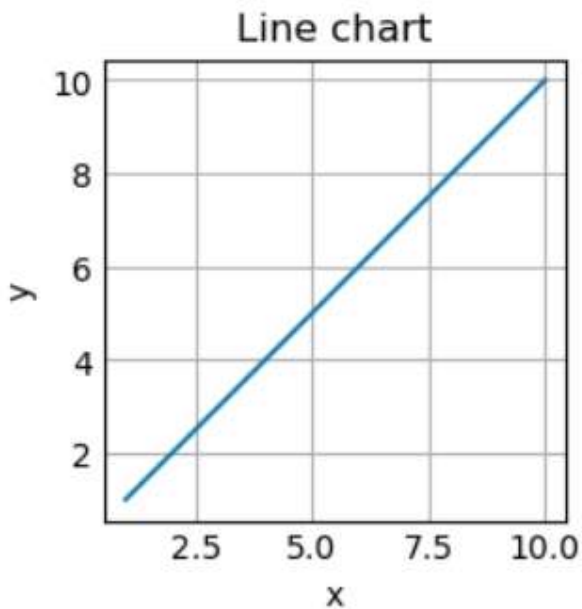
Line plots are the basic charts where dot consecutive points are connected by a continuous line.

The default plot() is used for line charts.

Code snippet:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 x=np.array([1,10])
5 y=np.array([1,10])
6
7 plt.plot(x,y)
8 plt.title("Line chart")
9 plt.xlabel("x")
10 plt.ylabel("y")
11 plt.show()
12
```

Output:



Description:

np.array() will generate an array in the specified range.

plot() for the plotting the line chart.

title() for defining the title, ylabel() and xlabel () for labelling y and x axis respectively.

show() displays the graph.

BAR GRAPH

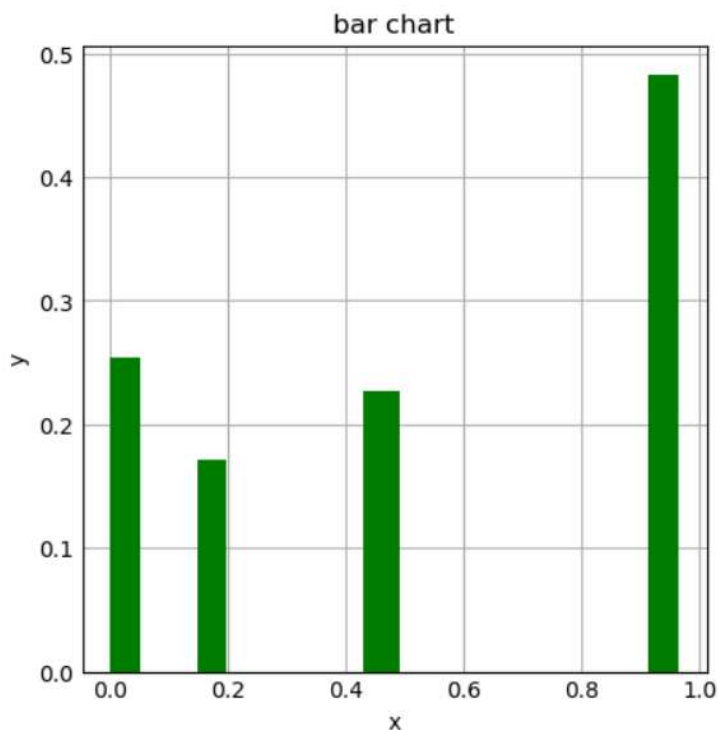
Rectangular bars are used in bar chart as the height represent the frequency of a particular element. bar() is used for plotting bar graphs, it has parameters like x, y, width, color.

bar(x, y, color, width)

Code snippet:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 x=np.random.rand(5)
5 y=np.random.rand(5)
6 print(x)
7
8 fig = plt.figure(figsize = (4, 4))
9 plt.bar(x,y, width=0.05, color="green")
10 plt.title("bar chart")
11 plt.xlabel("x")
12 plt.ylabel("y")
13 plt.show()
```

Output:



Description:

random.rand(5) generates a random array.

Here bar() is used along with the parameters, width which denoted the width of the rectangular bars and color for the bars.

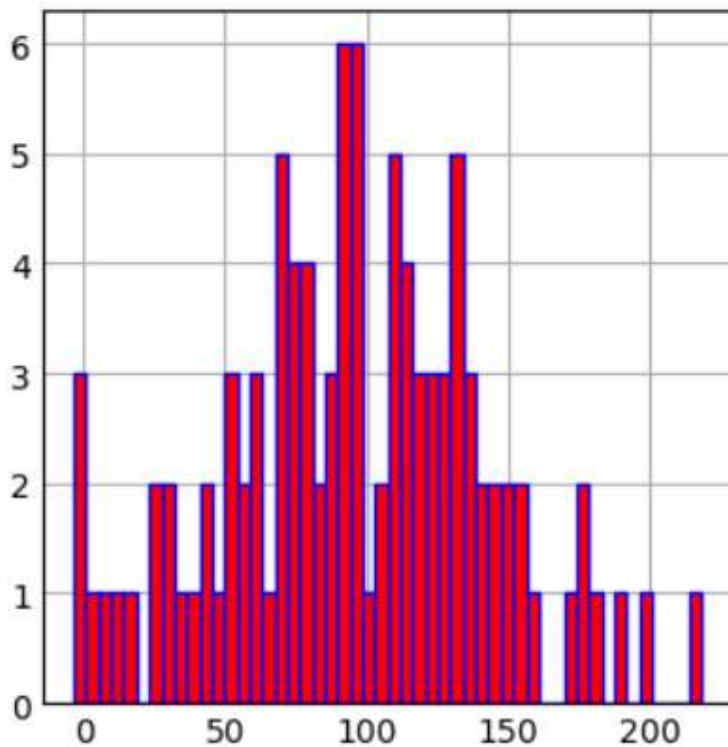
HISTOGRAM

Histogram is a type of bar graph where the graph is represented in groups. hist() is used for plotting histogram with parameters x, bins, color, edgecolor.

Code snippet:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 x=np.random.normal(100,50,100)
5 fig = plt.figure(figsize = (3, 3))
6 plt.hist(x,bins=50,color="red", edgecolor="blue")
7 plt.show()
```

Output:



Description:

hist() have parameters x which is a random generated array around one hundred with standard deviation 50 of 100 values.

bins indicate the number of sections on x axis.

color indicates the color of rectangular bars.

edge color that divides the rectangular bars.

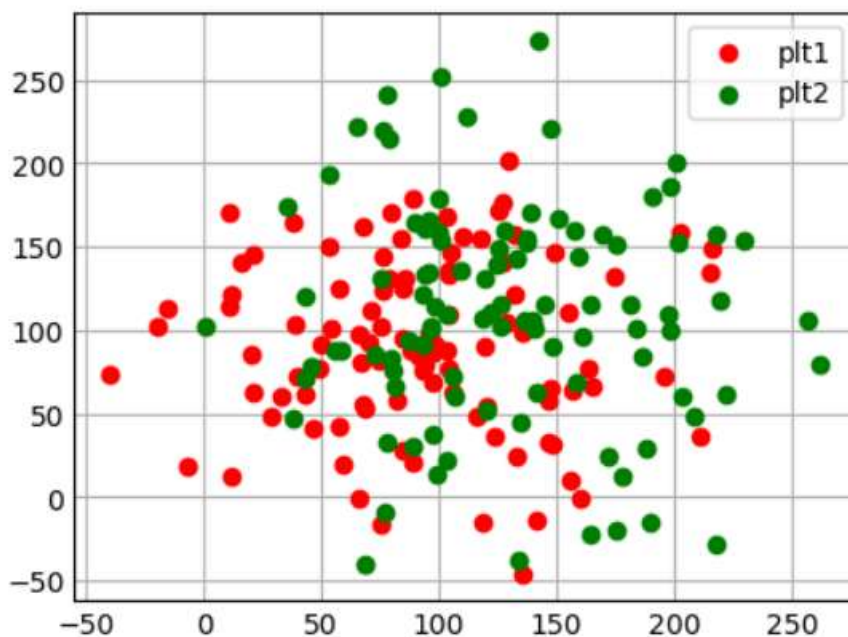
SCATTER PLOT

Scatter plot uses dots to depict the relationship between the data. `scatter()` is used for plotting scatter plot and scatter plot can be done for multiple datasets at the same time by differentiating it with colors. Legends help to achieve this difference.

Code snippet:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 x1=np.random.normal(100,50,100)
5 y1=np.random.normal(100,50,100)
6 x2=np.random.normal(120,60,100)
7 y2=np.random.normal(120,60,100)
8 fig = plt.figure(figsize = (4, 3))
9 plt.scatter(x1,y1,c="r")
10 plt.scatter(x2,y2,c="g")
11 plt.legend(["plt1","plt2"])
12
13 plt.show()
```

Output:



Description:

`x1, y1` belong to one dataset and `x2, y2` belong to another dataset.

Here the scatter() is used for the scatter plot and the parameters x, y and c which represents color.

Legend adds the label which helps to differentiate the plots.

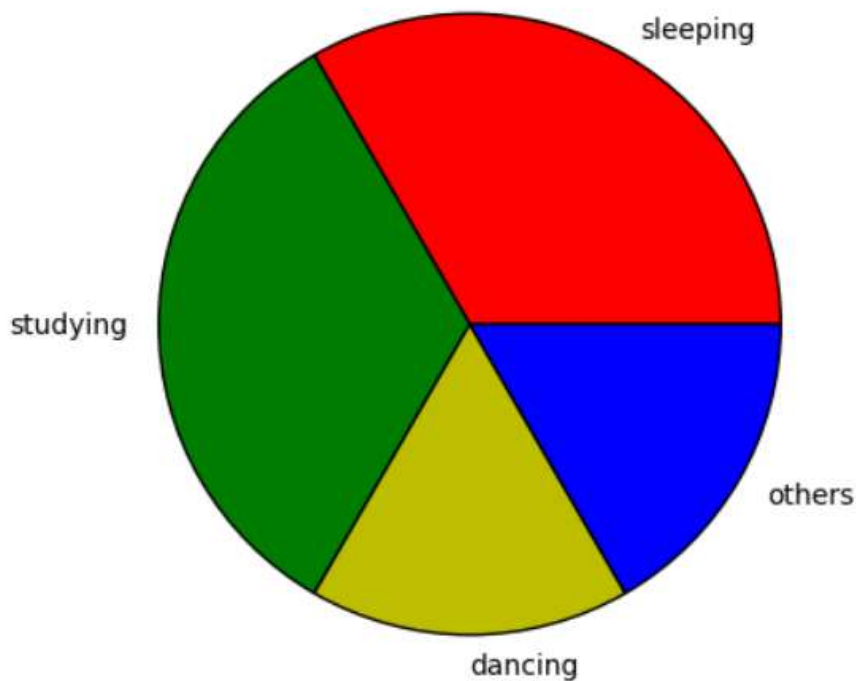
PIECHART

Pie charts are used to plot data of same kind which means the same series of data where the different elements are divide based on their percentage.

Code snippet:

```
1 from matplotlib import pyplot as plt
2 import numpy as np
3 time = [8,8,4,4]
4 color=["r","g","y","b"]
5 work= ["sleeping","studying","dancing","others"]
6 fig = plt.figure(figsize=(4, 4))
7 plt.pie(time, labels=work, wedgeprops={"edgecolor":"black","linewidth":1}, colors=color)
8 plt.show()
```

Output:



Description:

Here in the code, time is the x variable and labels are defined using labels keyword.

wedgeprops attribute is used to define the linewidth and the color to separate the elements in the pie chart.

color array contains the colors to each element.

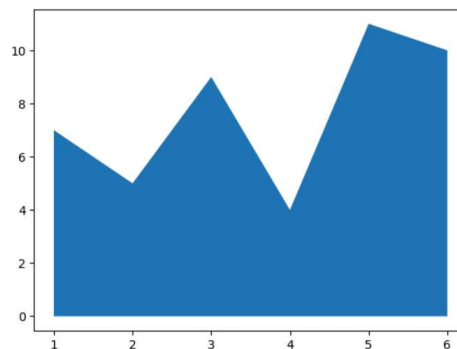
AREA CHART

In area chart, the area under the line to x axis is filled or shaded. It can be done by using `fill_between()` function or `stackplot()` function.

Code snippet:

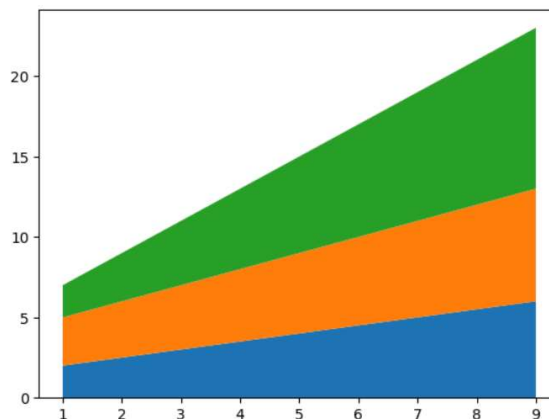
Using `fill_between()` function

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4
5 x=[1,2,3,4,5,6]
6 y=[7,5,9,4,11,10]
7
8 plt.fill_between(x, y)
9 plt.show()
```



Using `stackplot()` function

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 x=[1,3,5,7,9]
5 y1=[2,3,4,5,6]
6 y2=[3,4,5,6,7]
7 y3=[2,4,6,8,10]
8
9 plt.stackplot(x,y1, y2, y3)
10 plt.show()
```



Description:

`fill_between()` function have parameters of x and y.

`stackplot()` function parameters are x, y1,y1,y3 where y1, y2, y3 are the different groups.

Legends are added by mentioning the labels attribute in `stackplot()` function and giving the position to the legend using `plt.legend()` function with attribute loc.

SEABORN

Seaborn is a visualization library in python which is built on top of matplotlib. It is advanced than matplotlib and have different features which are default like style, color palette. Seaborn has different categories of plots like relational plot, categorical plot, distribution plot, regression plot, matrix plot.

Let us see the different plots in each category.

1 Relational Plots:

- scatterplot()
- lineplot()
- relplot()

2 Categorical Plots:

- barplot()
- countplot()
- boxplot()
- violinplot()
- swarmplot()
- pointplot()
- catplot()

3 Distribution Plots:

- histplot()
- kdeplot()
- rugplot()
- distplot()

4 Relational Plots:

- regplot()
- lmpplot()

5 Matrix Plots:

- heatmap()
- clustermap()

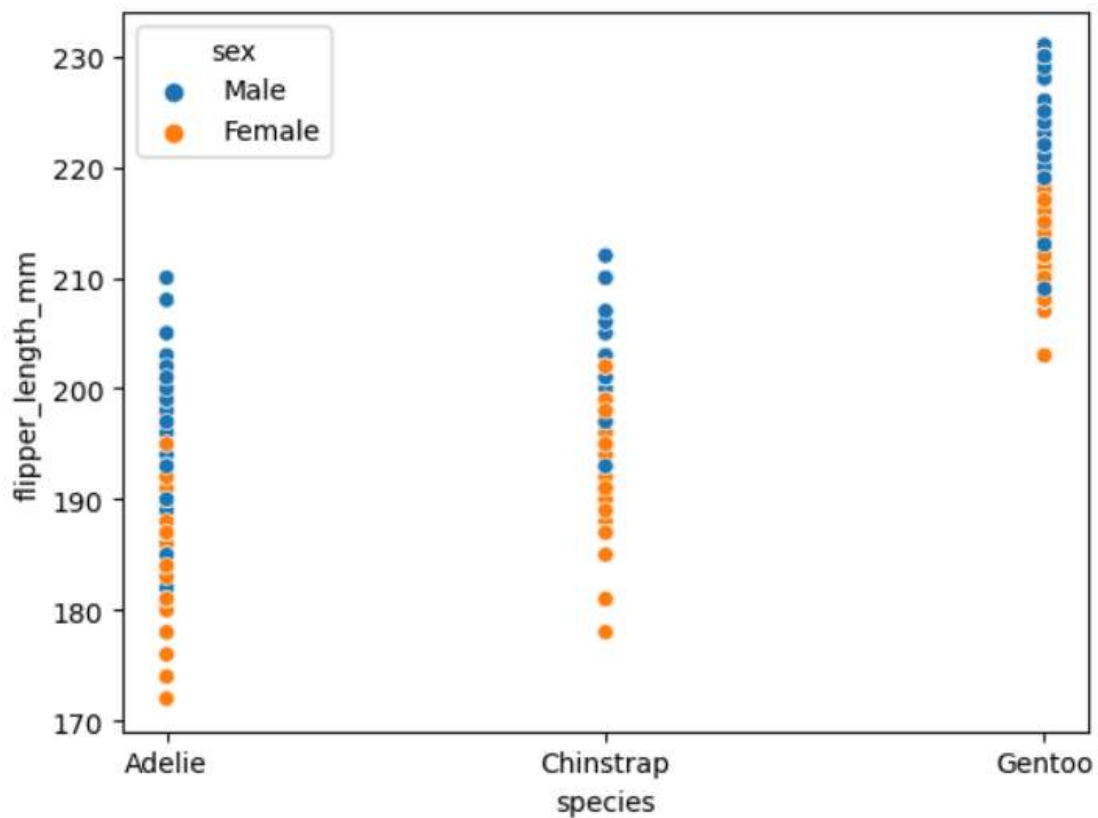
Here are some sample codes for some of the graphs.

Scatter plot:

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 df = sns.load_dataset("penguins")

1 sns.scatterplot(y='flipper_length_mm', x='species', hue='sex', data=df)
2 plt.show()
```

Output:



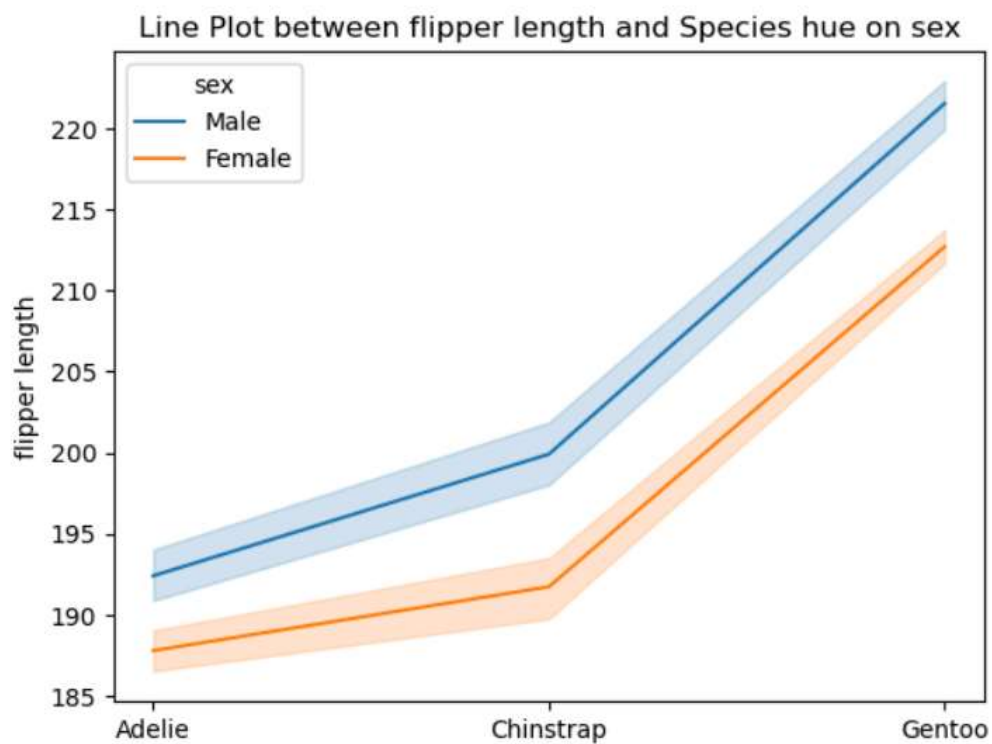
Description:

Here the different colors on sex is due to the parameter hue. Since hue is on “sex” column there is a difference for male and female.

LINE PLOT

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import seaborn as sns
4
5 sns.lineplot(y='flipper_length_mm', x='species', hue='sex', data=df)
6 plt.title('Line Plot between flipper length and Species hue on sex')
7 plt.ylabel('flipper length')
8 plt.xlabel('Species')
9 plt.show()
```

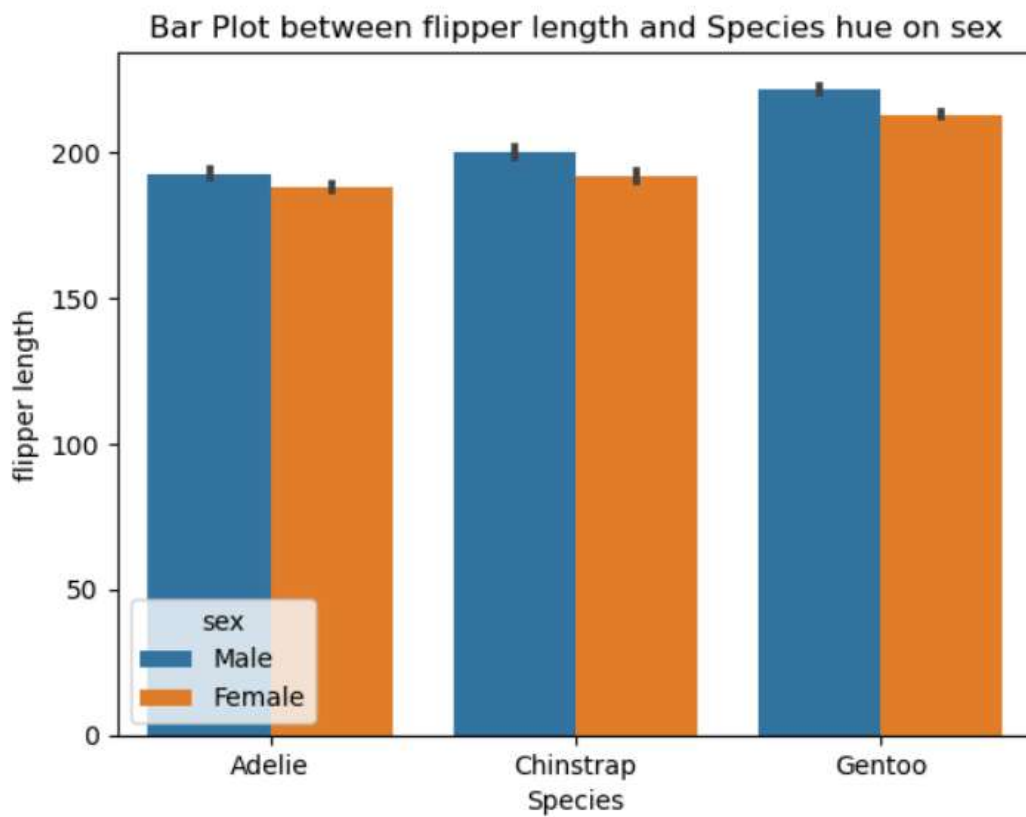
Output:



BAR PLOT

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import seaborn as sns
4
5 sns.barplot(y='flipper_length_mm', x='species', hue='sex', data=df)
6 plt.title('Bar Plot between flipper length and Species hue on sex')
7 plt.ylabel('flipper length')
8 plt.xlabel('Species')
9 plt.show()
```

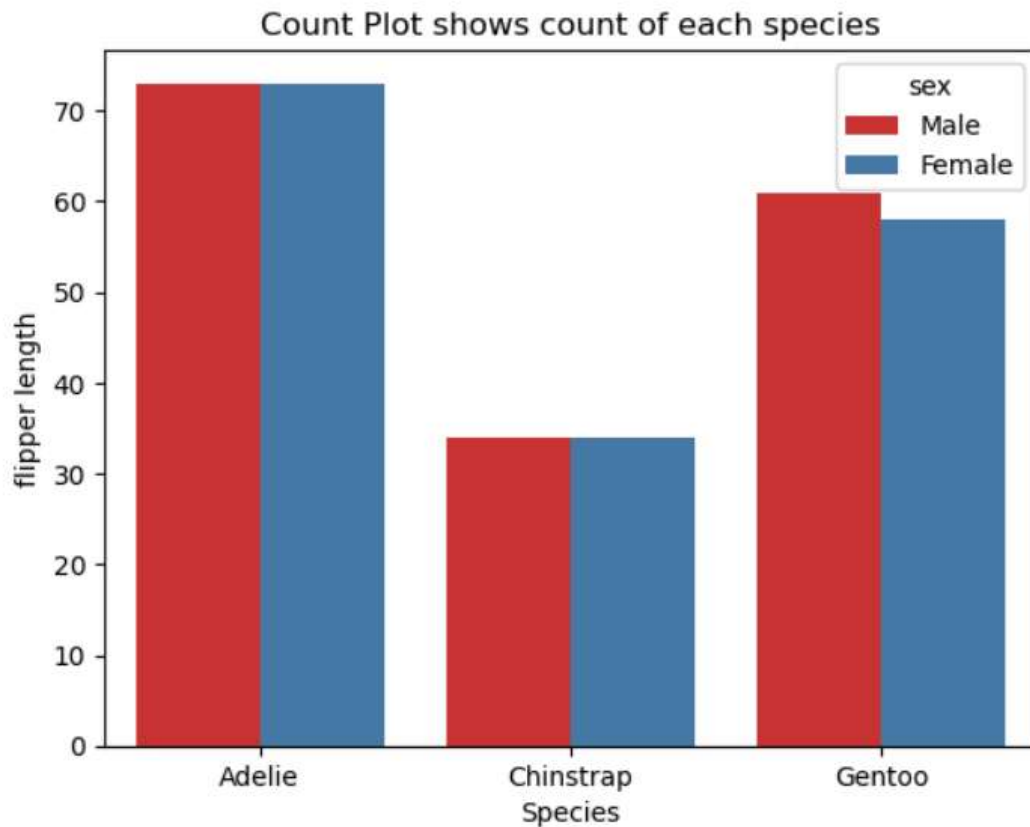
Output:



COUNT PLOT

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import seaborn as sns
4
5 sns.countplot(x='species', hue='sex', data=df, palette="Set1")
6 plt.title('Count Plot shows count of each species')
7 plt.ylabel('flipper length')
8 plt.xlabel('Species')
9 plt.show()
```

Output:



Description:

The count plots give the count of the data passed.

Here x is the species, hue is sex and palette as Set1.

Palette sets the color and it is Set1 is the default palette in seaborn.

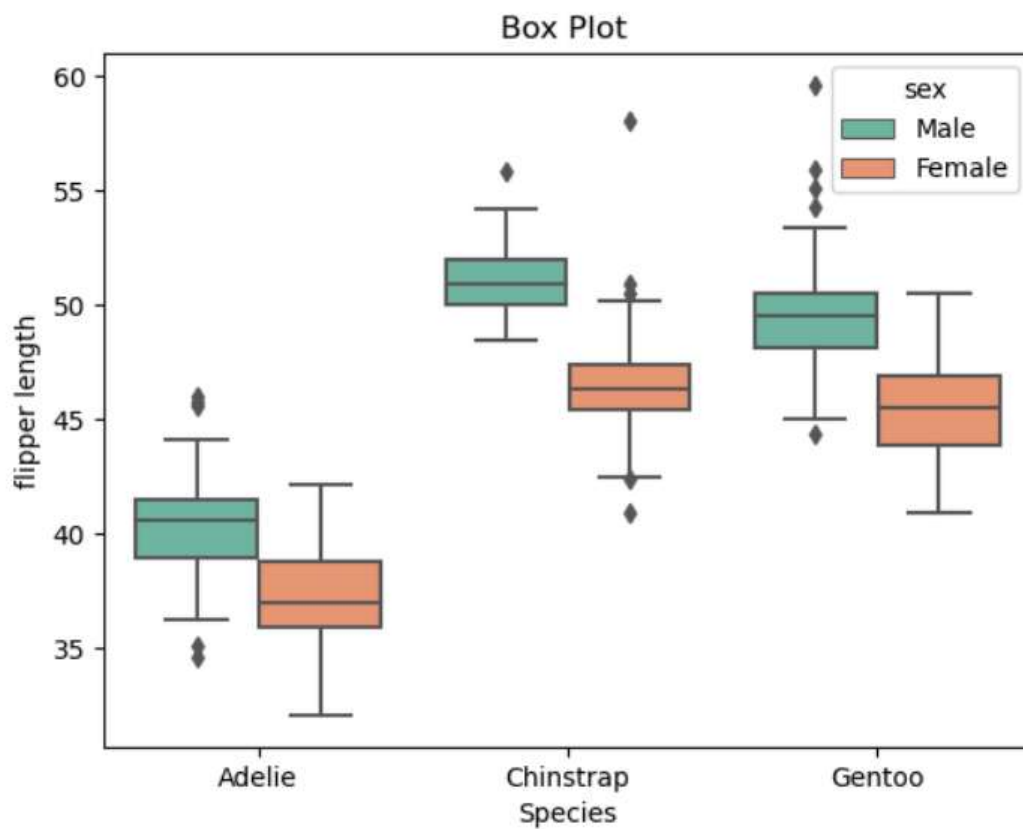
BOX PLOT

Box plot shows the quartiles where the data lies in interquartile range will be inside the box. The points which are away from the whiskers are outliers.

Code snippet:

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import seaborn as sns
4
5 sns.boxplot(x='species',y="bill_length_mm",hue='sex', data=df, palette="Set2")
6 plt.title('Box Plot ')
7 plt.ylabel('flipper length')
8 plt.xlabel('Species')
9 plt.show()
```

Output:



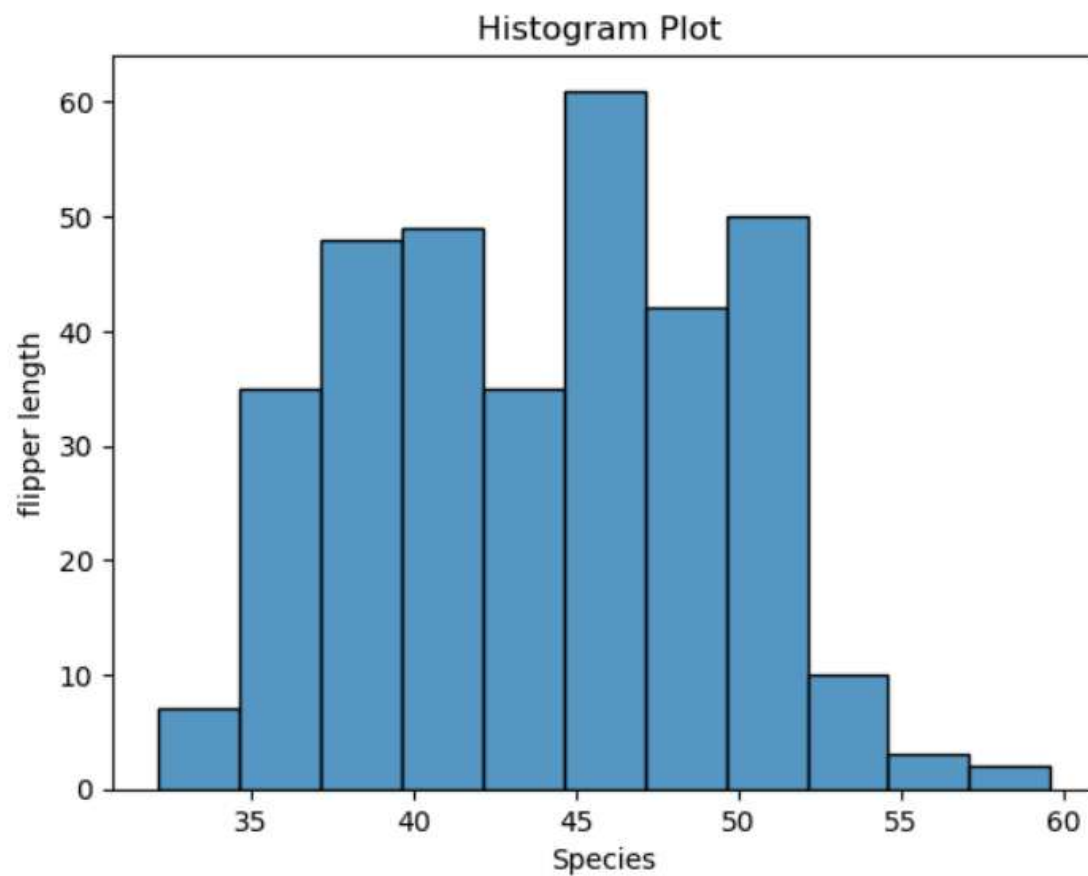
Description:

The parameter x should be numeric as the plot shows the difference between categories.

HISTOPLLOT

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import seaborn as sns
4
5 sns.histplot(x='bill_length_mm', data=df, palette="Set2")
6 plt.title('Histogram Plot ')
7 plt.ylabel('flipper length')
8 plt.xlabel('Species')
9 plt.show()
```

Output:



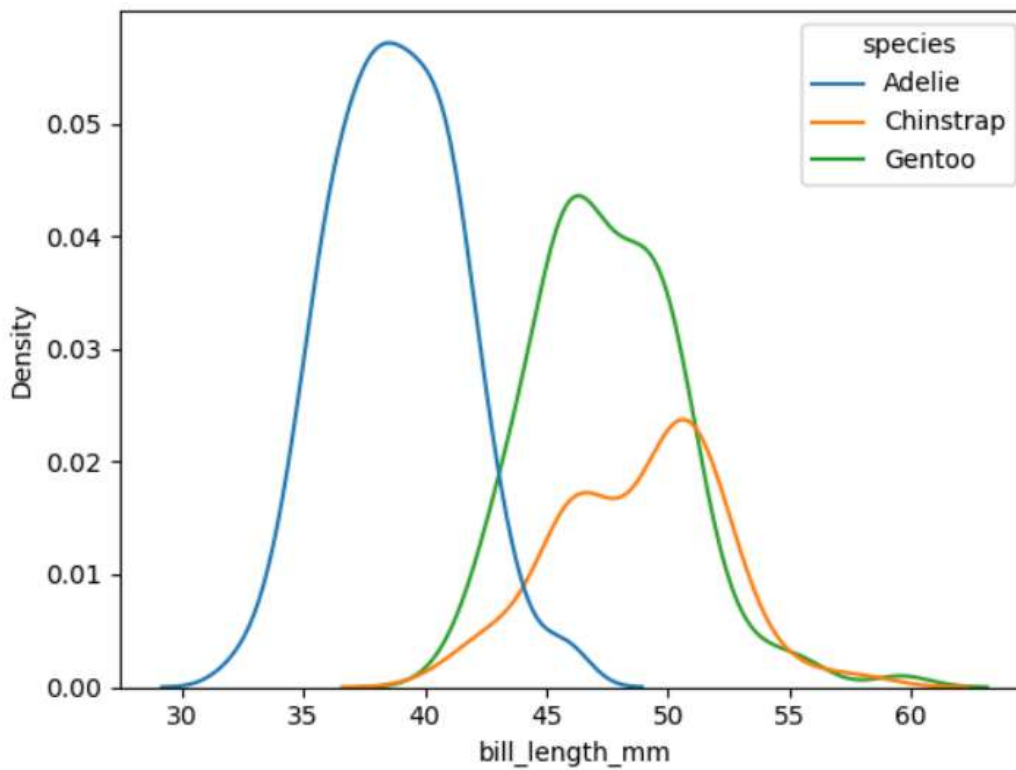
KDEPLOT

KDE stands for kernel density estimate. It creates a curve as based on probability. It creates single graph for multiple data samples.

Code snippet:

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import seaborn as sns
4
5 sns.kdeplot(x='bill_length_mm',hue='species',data=df )
6 plt.show()
```

Output:



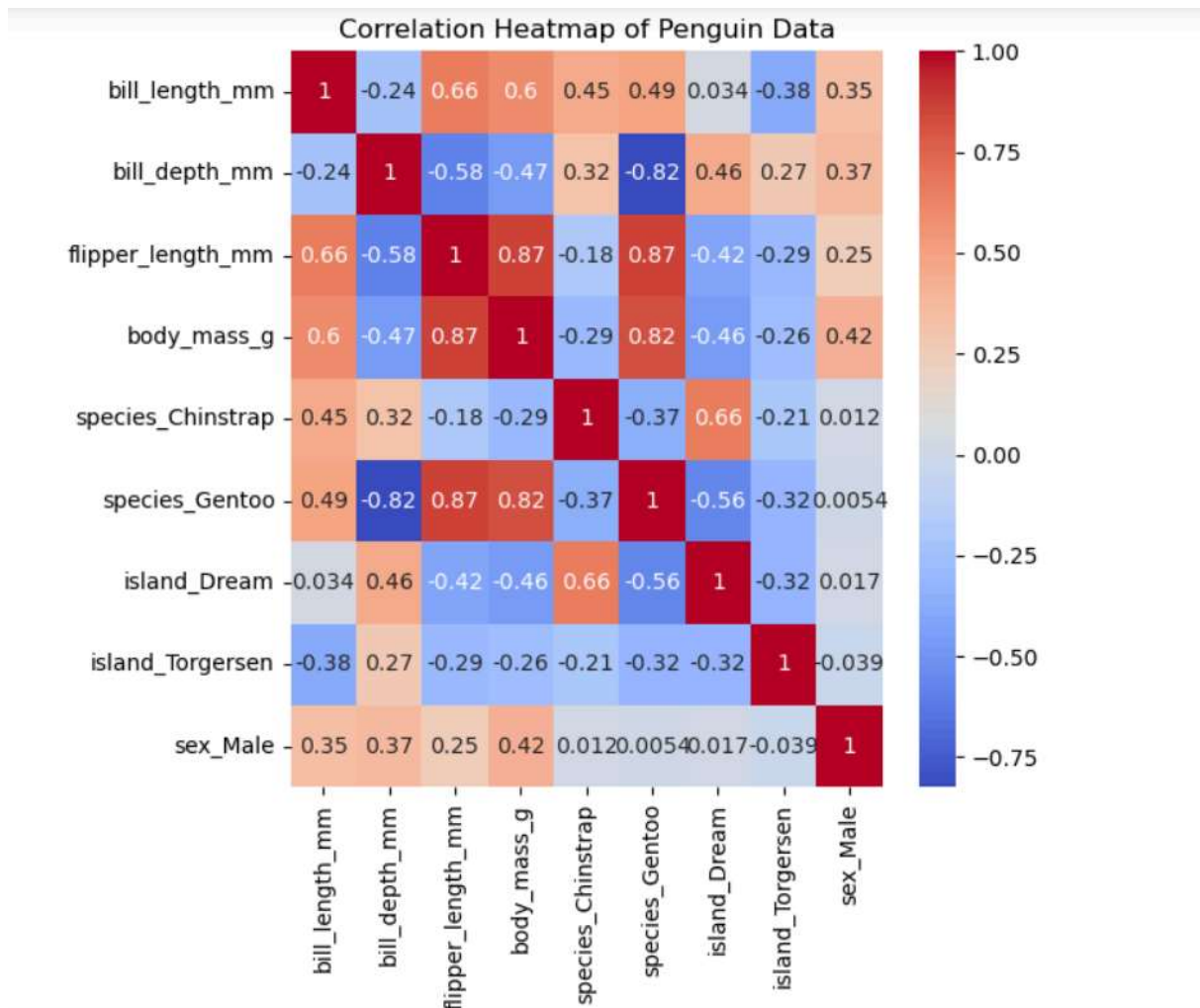
HEATMAP

Heatmaps uses correlation matrix and visualizes data. The datapoints where the higher values get brighter colors and lower values get darker colors.

Code snippet:

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 penguins_encoded = pd.get_dummies(df, columns=['species','island','sex'], drop_first=True)
6
7 correlation_matrix = penguins_encoded.corr()
8
9 plt.figure(figsize=(6, 6))
10 sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm")
11 plt.title("Correlation Heatmap of Penguin Data")
12 plt.show()
13
```

Output:



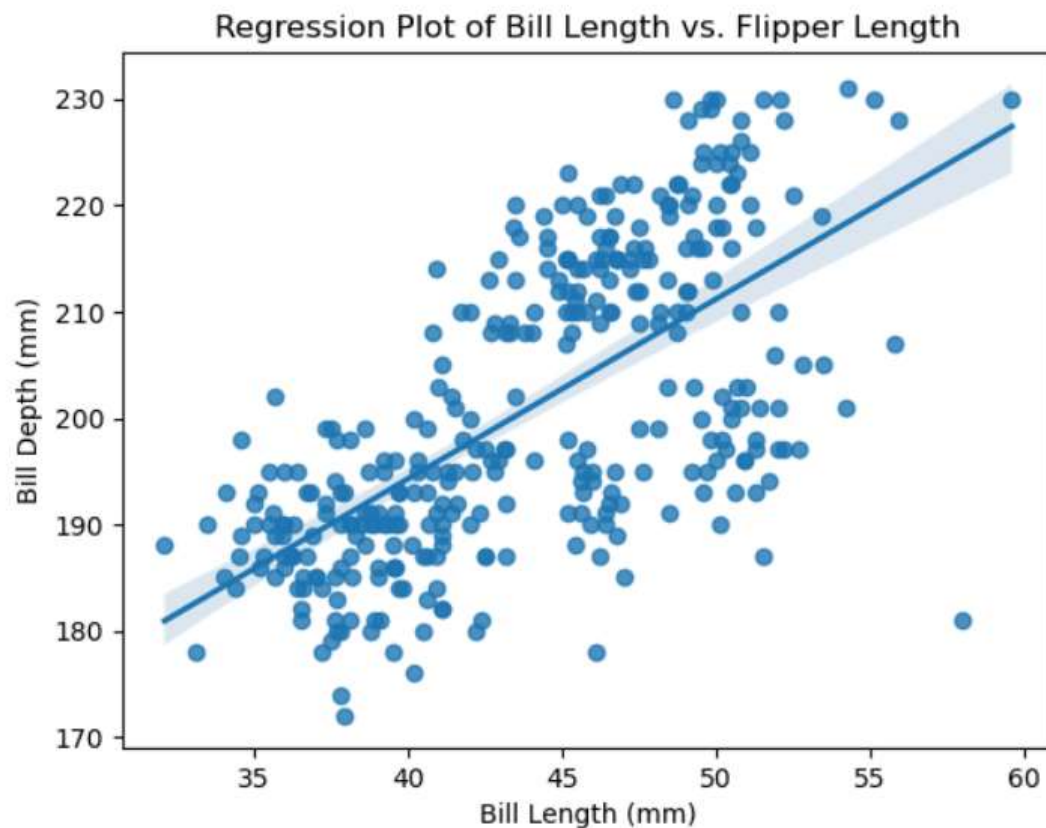
REGPLOT

Regression plots shows the relationship between variables along with the regression line.

Code snippet:

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 sns.regplot(x="bill_length_mm", y="flipper_length_mm", data=penguins)
5 plt.title("Regression Plot of Bill Length vs. Flipper Length")
6 plt.xlabel("Bill Length (mm)")
7 plt.ylabel("Bill Depth (mm)")
8 plt.show()
9
```

Output:



COMPARISON OF MATPLOTLIB AND SEABORN

Matplotlib:

- **Core Library:** Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python. It provides a low-level interface for creating plots with fine-grained control over every aspect of the figure.
- **Flexibility:** Matplotlib allows users to create any type of plot imaginable, from basic line plots and scatter plots to complex 3D plots and geographical maps.
- **Mature:** Matplotlib has been around for a long time and is a well-established library in the Python ecosystem. Many other libraries, including Seaborn, are built on top of Matplotlib.
- **Customization:** Matplotlib offers extensive customization options, allowing users to customize every aspect of a plot, including colors, line styles, markers, fonts, annotations, and more.
- **Integration:** Matplotlib can be easily integrated with other libraries and frameworks, such as NumPy, Pandas, SciPy, which makes it a versatile choice for data visualization in various contexts.

Advantages of Matplotlib:

- High degree of customization and control over plots.
- Wide range of plot types and styles.
- Strong community support and extensive documentation.
- Well-suited for creating publication-quality figures and graphics.

Seaborn:

- **High-Level Interface:** Seaborn is built on top of Matplotlib and provides a high-level interface for creating attractive and informative statistical graphics. It simplifies the process of creating complex visualizations by providing default aesthetics and convenient functions for common tasks.
- **Statistical Plotting:** Seaborn specializes in statistical plotting and offers a variety of plot types specifically designed for visualizing relationships in data, such as scatter plots, bar plots, box plots, violin plots, pair plots, and more.
- **Aesthetics:** Seaborn comes with built-in themes and color palettes that make it easy to create visually appealing plots with minimal effort. It also provides tools for fine-tuning plot aesthetics, such as controlling the size of plot elements and adjusting the color palette.
- **Integration with Pandas:** Seaborn seamlessly integrates with Pandas dataframes, allowing users to easily plot data stored in Pandas data structures without the need for extensive data manipulation.

Advantages of Seaborn:

- Simplified syntax and high-level functions for creating complex plots.
- Built-in support for statistical plotting and data exploration.
- Attractive default aesthetics and color palettes.
- Seamless integration with Pandas dataframes for data visualization.
- Ideal for exploratory data analysis and quick visualization of relationships in data.

Overall, the choice between Matplotlib and Seaborn depends on the specific requirements of data visualization task. Matplotlib is preferred for its flexibility and fine-grained control over plots, while Seaborn is preferred for its ease of use, statistical plotting capabilities, and attractive default aesthetics. Many users combine both libraries, leveraging the strengths of each as needed.