# Chapter 1

# Introduction to Seaborn

## 1.1 Definition and Overview

Seaborn is a sophisticated Python data visualization library built on top of Matplotlib. It is designed to create attractive and informative statistical graphics with simple, high-level commands. While Matplotlib provides the foundational building blocks for plotting, Seaborn offers a more streamlined and aesthetically pleasing interface, especially for common statistical tasks.

## 1.2 Why it matters

In the field of data science and analysis, effective visualization is crucial for exploratory data analysis (EDA), identifying patterns, communicating findings, and diagnosing models. Seaborn excels at this by integrating seamlessly with pandas DataFrames, allowing you to map data variables directly to plot aesthetics. It simplifies the creation of complex multi-plot grids and comes with beautiful default styles and color palettes.

## 1.3 Scope

Seaborn's functionality can be broadly categorized into visualizing distributions, categorical data, relationships between variables, and matrix data. This guide will provide a deep dive into each of these areas, focusing on the most common and powerful plotting functions.

# Chapter 2

# Deep Explanation and Examples

## 2.1 Distribution Plots

These plots are used to understand the distribution of a dataset, either for a single variable (univariate) or between two variables (bivariate).

### 2.1.1 distplot (and its successors displot and histplot)

**Deep Explanation:** The `distplot` function is a classic Seaborn tool for visualizing the univariate distribution of observations. It combines several plots into one: a histogram, a Kernel Density Estimate (KDE) plot, and optionally a rug plot.

- **Histogram:** Divides the data into bins and shows the frequency of data points in each bin.

- **KDE Plot:** Provides a smoothed, continuous line to estimate the probability density function of the variable.

   **Note:** In recent versions, `distplot` has been deprecated. Its functionality is now split between `displot` (a figure-level function) and `histplot` (an axes-level function), which offer more power and clarity. For modern code, it's recommended to use `histplot`.
   **Example:**

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

# Generate sample data
data = np.random.normal(size=100)

# Using histplot (the modern replacement for distplot)
sns.histplot(data, kde=True)
plt.title('Histogram with KDE for a Single Variable')
plt.xlabel('Value')
plt.ylabel('Frequency')
plt.show()
```

### 2.1.2   jointplot

**Deep Explanation:** This function creates a multi-panel figure that shows both the bivariate (or joint) relationship between two variables and the univariate (or marginal) distribution of each on separate axes. It's excellent for seeing a relationship and the individual distributions simultaneously. The central plot can be customized using the `kind` parameter ('scatter', 'reg', 'hex', 'kde').

**Example:**

```python
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt

# Load a sample dataset
tips = sns.load_dataset("tips")

# Create a jointplot to see the relationship between total_bill and tip
# 'kind="reg"' adds a scatter plot and a linear regression line
sns.jointplot(data=tips, x="total_bill", y="tip", kind="reg")
plt.suptitle('Joint Plot of Total Bill and Tip', y=1.02) # Adjust title
    position
plt.show()
```

### 2.1.3   pairplot

**Deep Explanation:** The `pairplot` function is a powerful tool for visualizing pairwise relationships in a dataset. It creates a grid of axes such that each numerical variable in the data is shared across the y-axes on a single row and the x-axes on a single column. The diagonal plots show the univariate distribution of each variable (as a histogram or KDE), while the off-diagonal plots show the bivariate relationship (as a scatter plot). The `hue` parameter can be used to color the data points based on a categorical variable.

**Example:**

```python
import seaborn as sns
import matplotlib.pyplot as plt

# Load the iris dataset
iris = sns.load_dataset("iris")

# Create a pairplot to see relationships between all numerical features
    , colored by species
sns.pairplot(iris, hue="species", palette="viridis")
plt.suptitle('Pairwise Relationships in the Iris Dataset', y=1.02)
plt.show()
```

### 2.1.4   rugplot

**Deep Explanation:** A `rugplot` is a very simple plot that draws a small vertical tick at each observation along an axis. It is used to visualize the distribution of a single variable. While it can be used on its own, it is more effective when used as a component of a more complex plot (like `distplot`) to show the location of individual observations.

**Example:**

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 # Sample data
6 data = np.random.randn(50)
7
8 # Create a rugplot
9 sns.rugplot(x=data, height=0.5)
10 plt.title('Rug Plot Showing Individual Data Points')
11 plt.show()
```

### 2.1.5  kdeplot

**Deep Explanation:** The `kdeplot` (Kernel Density Estimate plot) is used to visualize the probability density of a continuous variable. It can be thought of as a smoothed version of a histogram. It can be used for a single variable (1D) or to show the joint probability density of two variables (2D). A 2D KDE plot uses contours or shading to represent areas of higher density.

**Example:**

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Load the tips dataset
5 tips = sns.load_dataset("tips")
6
7 # Create a 2D KDE plot for total_bill and tip
8 sns.kdeplot(data=tips, x="total_bill", y="tip", fill=True, cmap="mako")
9 plt.title('2D KDE Plot of Total Bill and Tip')
10 plt.show()
```

## 2.2  Categorical Data Plots

These plots help in visualizing the relationship between a numerical variable and one or more categorical variables.

### 2.2.1  factorplot (now catplot)

**Deep Explanation:** `factorplot` has been renamed to `catplot` in recent versions of Seaborn. It is a figure-level function that provides a unified interface for several axes-level categorical plots. Its primary strength is its ability to create "faceted" plots across different subsets of the data using the `row` and `col` parameters. You select the specific type of plot using the `kind` parameter (e.g., 'box', 'violin', 'bar', 'strip').

**Example:**

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Load the tips dataset
5 tips = sns.load_dataset("tips")
```

```
6
7  # Using catplot (the modern factorplot) to create faceted boxplots
8  # kind='box' specifies a boxplot
9  # col='time' creates separate columns for 'Lunch' and 'Dinner'
10 sns.catplot(data=tips, x="day", y="total_bill", kind="box", col="time")
11 plt.suptitle('Total Bill by Day, Faceted by Time of Day', y=1.02)
12 plt.show()
```

### 2.2.2   boxplot

**Deep Explanation:** A box plot (or box-and-whisker plot) provides a five-number summary of a set of data: the minimum, first quartile (Q1), median (Q2), third quartile (Q3), and maximum. The "box" represents the interquartile range (IQR), which contains the middle 50% of the data. The "whiskers" typically extend to 1.5 times the IQR, and any points beyond them are considered outliers. It's excellent for comparing distributions across different categories.

**Example:**

```
1  import seaborn as sns
2  import matplotlib.pyplot as plt
3
4  # Load the tips dataset
5  tips = sns.load_dataset("tips")
6
7  # Create a boxplot to compare total_bill distributions across different
       days
8  sns.boxplot(data=tips, x="day", y="total_bill", hue="smoker")
9  plt.title('Total Bill Distribution by Day and Smoker Status')
10 plt.show()
```

### 2.2.3   violinplot

**Deep Explanation:** A `violinplot` is a more advanced plot that combines a boxplot with a KDE plot. The "violin" shape shows the probability density of the data at different values. This provides more information than a standard boxplot, as it can show if a distribution is bimodal (has two peaks) or skewed, which a boxplot might hide.

**Example:**

```
1  import seaborn as sns
2  import matplotlib.pyplot as plt
3
4  # Load the tips dataset
5  tips = sns.load_dataset("tips")
6
7  # Create a violin plot
8  # 'split=True' allows comparison of 'smoker' within the same violin
9  sns.violinplot(data=tips, x="day", y="total_bill", hue="smoker", split=
      True)
10 plt.title('Violin Plot of Total Bill by Day and Smoker Status')
11 plt.show()
```

## 2.2.4 stripplot

**Deep Explanation:** A `stripplot` is essentially a scatter plot where one of the variables is categorical. It shows every single observation. To prevent points from overlapping and obscuring the distribution, it's common to add a small amount of random "jitter" along the categorical axis. This is a good choice for smaller datasets.

**Example:**

```python
import seaborn as sns
import matplotlib.pyplot as plt

# Load the tips dataset
tips = sns.load_dataset("tips")

# Create a stripplot with jitter
sns.stripplot(data=tips, x="day", y="total_bill", jitter=True, hue="sex")
plt.title('Strip Plot of Total Bill by Day')
plt.show()
```

## 2.2.5 swarmplot

**Deep Explanation:** A `swarmplot` is similar to a `stripplot` but adjusts the points along the categorical axis with a non-overlapping algorithm. This gives a better representation of the distribution of values, especially for small to medium-sized datasets, without the randomness of jitter. However, it does not scale well to very large datasets.

**Example:**

```python
import seaborn as sns
import matplotlib.pyplot as plt

# Load the tips dataset
tips = sns.load_dataset("tips")

# Create a swarmplot
sns.swarmplot(data=tips, x="day", y="total_bill", hue="sex", dodge=True)
plt.title('Swarm Plot of Total Bill by Day')
plt.show()
```

## 2.2.6 barplot

**Deep Explanation:** A `barplot` shows an estimate of central tendency (by default, the mean) for a numeric variable, with the height of each bar representing the estimate. It also shows an indication of uncertainty around that estimate using error bars (which represent the confidence interval). It is useful for comparing an aggregate value across different categories.

**Example:**

```python
import seaborn as sns
import matplotlib.pyplot as plt

# Load the titanic dataset
titanic = sns.load_dataset("titanic")
```

```
6
7  # Create a barplot showing the average age for each passenger class
8  sns.barplot(data=titanic, x="class", y="age")
9  plt.title('Average Age by Passenger Class')
10 plt.show()
```

### 2.2.7   countplot

**Deep Explanation:** A `countplot` is essentially a histogram for a categorical variable. Instead of showing a statistic like the mean, it shows the total count of observations in each category. It's a simple yet effective way to see the frequency distribution of categorical data.

   **Example:**

```
1  import seaborn as sns
2  import matplotlib.pyplot as plt
3
4  # Load the titanic dataset
5  titanic = sns.load_dataset("titanic")
6
7  # Create a countplot to see the number of passengers in each class
8  sns.countplot(data=titanic, x="class", hue="who")
9  plt.title('Number of Passengers by Class')
10 plt.show()
```

## 2.3   Matrix Plots

Matrix plots are a way of visualizing data as a color-encoded matrix, where different colors represent different values.

### 2.3.1   heatmap

**Deep Explanation:** A `heatmap` is a graphical representation of data where individual values in a matrix are represented as colors. It is an incredibly useful tool for visualizing correlation matrices, confusion matrices, or any other tabular data. The color intensity indicates the magnitude of the value, making it easy to spot patterns, clusters, and high-or-low-value areas at a glance.

   **Example:**

```
1  import seaborn as sns
2  import matplotlib.pyplot as plt
3  import pandas as pd
4
5  # Load the flights dataset and pivot it into a matrix format
6  flights = sns.load_dataset("flights")
7  flights_pivot = flights.pivot(index="month", columns="year", values="
      passengers")
8
9  # Create a heatmap
10 plt.figure(figsize=(10, 8))
11 sns.heatmap(flights_pivot, cmap="magma", linecolor='white', linewidths
      =1, annot=True, fmt='d')
```

```
12 plt.title('Heatmap of Airline Passengers by Month and Year')
13 plt.show()
```

# Chapter 3

# Additional Resources and Summary

## 3.1  Related Concepts

- **Matplotlib:** The foundation of Seaborn. Knowing Matplotlib allows for deep customization of Seaborn plots.

- **Pandas:** Seaborn is designed to work with Pandas DataFrames. Proficiency in Pandas data manipulation is essential.

- **Statistical Aggregation:** Concepts like mean, median, standard deviation, and confidence intervals are integral to understanding what plots like `barplot` and `boxplot` represent.

- **Tidy Data:** Seaborn works best with "tidy" data, where each row is an observation and each column is a variable.

## 3.2  Assignments / Practice Questions

- **MCQ:** Which plot would you use to show the distribution of a single numerical variable while also displaying the underlying individual data points as ticks on the axis?

    1. boxplot
    2. histplot with rug=True
    3. barplot
    4. heatmap

- **Short Question:** Explain the key difference between a `barplot` and a `countplot`. When would you use one over the other?

- **Problem-Solving:** Using the titanic dataset (`sns.load_dataset("titanic")`), create a `violinplot` that shows the distribution of passenger age for each class ('First', 'Second', 'Third'). Further, split the violins to compare sex within each class.

- **Problem-Solving:** Create a correlation matrix for the numerical columns in the iris dataset (`sns.load_dataset("iris")`). Then, visualize this matrix using a `heatmap`. Annotate the heatmap with the correlation values.

- **Case Study:** You are analyzing a dataset of student test scores. The data includes the score, the subject ('Math', 'Science', 'History'), and the `study_hours`. Your goal is to see how score is related to `study_hours` but also see if this relationship differs by subject. Which plot would you use and why? Provide the code to generate it. (Hint: A figure-level function with faceting might be best).

## 3.3 Applications

- **Exploratory Data Analysis (EDA):** Seaborn is a primary tool for quickly understanding data distributions, correlations, and relationships before formal modeling.

- **Business Intelligence:** Creating insightful dashboards that visualize key performance indicators (KPIs) across different categories (e.g., sales by region).

- **Financial Analysis:** Visualizing stock correlations, portfolio performance, and economic indicators.

- **Scientific Research:** Publishing publication-quality plots of experimental results, comparing distributions between control and experimental groups.

- **Machine Learning:** Analyzing feature distributions, feature importance, and model performance (e.g., using a heatmap for a confusion matrix).

## 3.4 Related Study Resources

- Official Seaborn Tutorial: https://seaborn.pydata.org/tutorial.html

- Seaborn API Reference: https://seaborn.pydata.org/api.html

- Python Data Science Handbook by Jake VanderPlas (Chapter on Visualization with Seaborn): https://jakevdp.github.io/PythonDataScienceHandbook/04.14-visualization-with-seaborn.html

- Coursera - Applied Plotting, Charting & Data Representation in Python: https://www.coursera.org/learn/python-plotting

## 3.5 Summary / Key Takeaways

- **Foundation:** Seaborn is a statistical plotting library built on Matplotlib that works best with Pandas DataFrames.

- **Plot Categories:**

  - **Distribution Plots** (`histplot`, `jointplot`, `pairplot`): For visualizing the distribution and relationships of numerical data.

- **Categorical Plots** (`catplot`, `boxplot`, `barplot`, `swarmplot`): For comparing numerical data across different categorical groups.

- **Matrix Plots** (`heatmap`): For visualizing matrix-like data, especially correlations.

- **Key Strength:** Simplifies the creation of complex, multi-plot statistical graphics with concise syntax.

- **Figure-level vs. Axes-level:** Functions like `catplot` and `jointplot` are "figure-level" (they create a whole figure with potentially multiple subplots), while functions like `boxplot` and `scatterplot` are "axes-level" (they draw on a specific subplot).