

## PLOTTING WITH SEABORN:

SEABORN: It is a python data visualization library based on matplotlib.

- It provides a high-level interface for drawing attractive and informative statistical graphics.

### 1. UNIVARIATE ANALYSIS (FOR NUMERICAL VARIABLES)

i.e., Analyzing one variable at a time.

```
import numpy as np
```

```
import pandas as pd
```

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

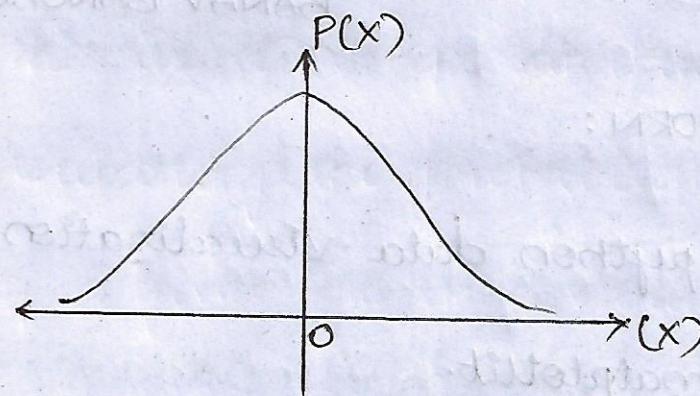
```
import seaborn as sns
```

→ importing seaborn from the

library.

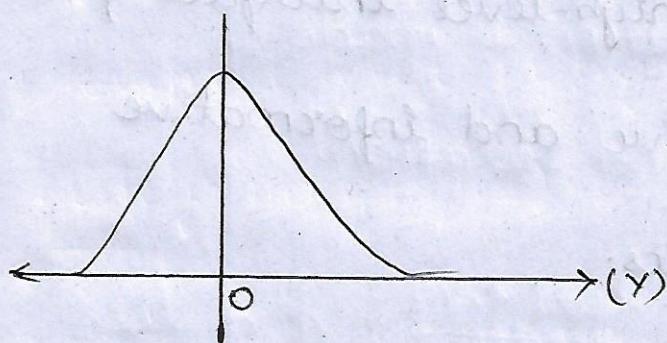
## EXAMPLE OF A NORMAL DISTRIBUTION:

(22)



$$\rightarrow \mu_x = 0$$

$$\sigma_x = 1$$



$$\rightarrow \mu_y = 1$$

$$\sigma_y = 1$$

The above graphs are symmetric.

```
1/p - sns.distplot(df['Y'])
```

```
# by default sns.distplot(x, hist=True, kde=True)
```

HISTOGRAM:

A histogram is a plot that lets us discover and show the underlying frequency distribution (shape) of a set of continuous data.

(23)

↳ This allows the inspection of the data for its underlying distribution (normal distribution), outliers, skewness ... etc.

NOTE :

\* IF BIN WIDTH IS TOO SMALL :

It shows too much individual data and does not allow the underlying pattern (frequency distribution) of the data.

\* IF BIN WIDTH IS TOO LARGE :

We are unable to find the underlying (pattern) trend in the data.

I/p - `sns.distplot(df['y'], kde=False, rug=True)`

# This is basically a frequency plot.

# Try if `rug=False` (`hist=True` is taken automatically as a default)

## KDE : KERNEL DENSITY ESTIMATION : (24)

This is a non-parametric method for estimating the probability density function of a given random variable.

- ↳ In general, they are generalisation and improvement over histograms (smoothering of histogram graphs)
- ↳ Here we all extract the important features of the data.

`rug` : bool, optional (whether to draw a rugplot on the support axis)

A rug plot is a plot of data for a single quantitative variable displayed as marks along an axis.

- ↳ As such it is an analogous to a histogram with zero-width bins, or a one-dimensional scatter plot.

(25)  $\text{Yp} = \text{sns.distplot}(\text{df}['y'], \text{hist=False}, \text{rug=True})$

OR

$\text{Yp} = \text{sns.distplot}$

UNDERSMOOTHED - Bandwidth is too small

OVERSMOOTHED - Bandwidth is too large

## 2. BIVARIATE ANALYSIS (FOR NUMERICAL VARIABLES)

$\text{Yp} = \text{df.describe()}$

$$\mu_x = 0, \mu_y = 0$$

$$\sigma_x = 1, \sigma_y = 1$$

Since the mean and median are nearly equal to '0' and '1', then there are no outliers.

### PROPERTIES OF KDE:

↳ Smooth

↳ No end points

↳ depends on bandwidth

### PROPERTIES OF HISTOGRAM:

↳ Not smooth

↳ depends on width of bins

↳ depends on endpoints of bins

## UNIVARIATE (Numerical) :

(26)

→ Histogram

→ KDE

→ Box

## BIVARIATE (Numerical variables) :

→ scatter

→ Hexbin

→ Pair plot.

### SCATTER PLOT:

These are similar plot to line graphs in that they use horizontal and vertical axes to plot data points.

→ This shows how much one variable is affected by another. This relationship between two variables is called their correlation.

→ Scatter plots usually consist of a large body of data.

(27) ↳ The closer the data points come when plotted to making a straight line.

↳ The higher the correlation between the two variables or the stronger the relationship.

↳ If the data points make a straight line going from the origin out to high x- and y-values, then the variables are said to have a positive correlation.

↳ If the line goes from a high-value on the y-axis down to a high-value on the x-axis, then the variables have a negative correlation.

i/p - `sns.jointplot(x='x', y='y', data=df, kind='scatter')`

or

`sns.jointplot(data=df, x='x', y='y', kind='scatter')`

## HEXBIN PLOTS: TO KNOW THE DENSITY (28)

A Hexbin plot is useful to represent the relationship of 2 numerical variables when we have a lot of data point.

↳ Instead of overlapping, the plotting window is split in several hexbins, and the number of points per hexbin is counted.

↳ The color denotes ~~this~~ the number of points.

NOTE:

\* IF LIGHTER : less Density

\* IF DARKER : High Density - most of the points lie in this region.

```
!/p - sns.jointplot(x='x', y='y', data=df, kind='hex',  
color='k')
```

or

```
sns.jointplot(data=df, x='x', y='y', kind='hex', color='k')
```

## PAIR PLOT:

(29)

A pair plot allows us to see both distribution of single variables and relationships between two variables.

→ It is a great method to identify trends for follow-up analysis and fortunately are easily implemented.

I/P - `iris = pd.read_csv('data/Iris.csv')`

`sns.pairplot(iris)`

O/P - Shows all the possible pairs plots.  
of

NOTE : Except species column, the rest columns plots are drawn.

I/P - `sns.pairplot(iris, hue = 'Species')`

hue - COLOR CODE and works with categorical data.

## WORKING WITH CATEGORICAL DATA: (30)

I/p - tips = pd.read\_csv('data/tips.csv')

tips.head()

I/p - tips.shape() → O/p - (244, 7)

I/p - tips.describe(include = 'all')

I/p - tips['day'].value\_counts()

O/p - Sat 87

Sun 76

Thur 62

Fri 19

Name: day, dtype: int64

### STRIP PLOT:

A strip plot is a scatter plot where one of the variables is categorical.

For example,

A box plot with an overlaid strip plot becomes more similar to a violin plot because some additional information about

(31) how the underlying data is distributed becomes visible.

I/P - `sns.stripplot(data=tips, x='data', y='total_bill')`

### SWARM PLOT :

This function is similar to `stripplot()`, but the points are adjusted (only along the categorical axis) so that they don't overlap.

→ It gives a better representation of the distribution of values, but it does not scale well to large numbers of observations.

→ This style of plot is sometimes called as a "BEE SWARM".

I/P - `sns.swarmplot(data=tips, x='day', y='total_bill')`

## BOXPLOT:

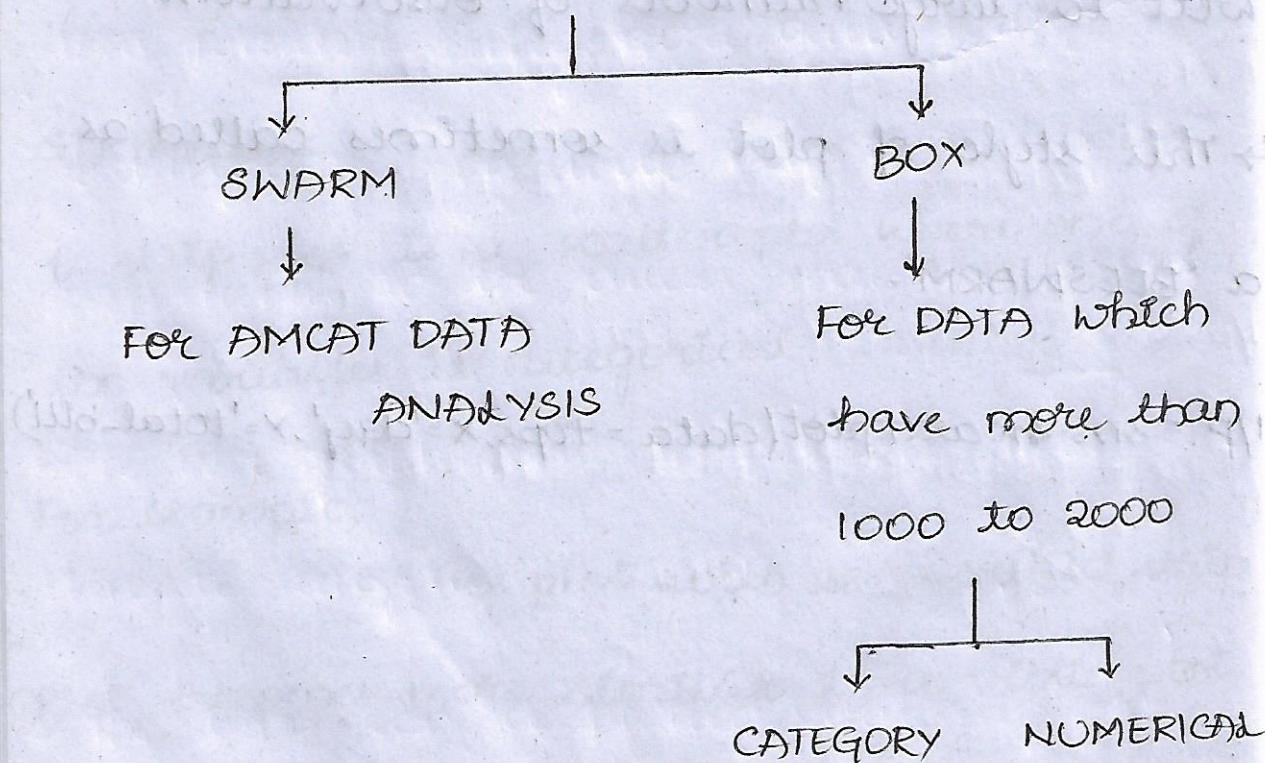
(32)

A box plot is a method for graphically depicting groups of numerical data through their quartiles.

↳ The box extends from the Q<sub>1</sub> to Q<sub>3</sub> quartile values of the data, with a line at the median (Q<sub>2</sub>).

↳ The whiskers extend from the edges of box to show the range of the data.

## USES



Q3  
I/p - sns.boxplot(data=tips, x='day', y='total\_bill')

I/p - sns.tips.loc[tips.time == 'Lunch']

In SQL: Select \* (all) from tips where time =  
= lunch

I/p - tips.loc[tips.time == 'lunch', 'day'].value-  
counts()

O/p - There 61  
For 7

Name: day, dtype: int64

In SQL: select \*<sup>day</sup> from tips where time =  
= lunch.

I/p - sns.boxplot(data=tips, x='day', y='total\_bill',  
hue='time')

# Adding a new column, indicating  
weekend or not.

I/p - tips['weekend'] = tips['day'].isin(['sat', 'sun'])  
tips.head()

O/P -

(34)

	total_bill	tip	sex	smoker	day	time	size	Weekend
0	16.99	1.01	F	No	Sun	Dinner	2	True
1	10.34	1.66	M	No	Sun	Dinner	3	True
2	21.01	3.50	M	No	Sun	Dinner	3	True
3	23.68	3.31	M	No	Sun	Dinner	2	True
4	24.59	3.61	F	No	Sun	Dinner	4	True

I/P - sns.boxplot(data=tips, x='day', y='total\_bill', hue='weekend')

#### BAR PLOT :

A bar chart or graph is a chart or graph that presents categorical data with rectangular bars with heights or lengths proportional to the values that they represent.

→ The bar graph shows comparisons among the discrete categories.

→ The bars can be plotted vertically (or) horizontally.

! / p - titanic = pd.read\_csv('data/titanic.csv')

titanic.head()

! / p - titanic.shape

o/p - (891, 15)

! / p - titanic.describe(include = 'all')

! / p - sns.barplot(x = 'sex', y = 'age', data = titanic)

Barplot → Estimator - mean (by default)

↳ Estimator - median

! / p - plt.figure(figsize=(12, 6))

plt.subplot(1, 2, 1)

sns.barplot(x = 'sex', y = 'age', data = titanic)

plt.title('Avg Age')

plt.subplot(1, 2, 2)

sns.barplot(x = 'sex', y = 'age', data = titanic,  
estimator=np.median)

plt.title('Median Age')

plt.show()

```
!p - sns.boxplot(x='sex', y='age', data=titanic) (36)
```

```
!p - sns.barplot(x='sex', y='survived', data=titanic)
```

```
!p - sns.barplot(x='sex', y='survived', hue='class',  
data=titanic)
```

### COUNT PLOT:

A count plot can be thought of as a histogram across a categorical, instead of quantitative, variable.

↳ The basic API and options are identical to those for barplot(), so you can compare counts across nested variables.

```
!p - sns.countplot(x='class', data=titanic,  
palette='Greens_d')
```

### PALETTE = COLOR REQUIREMENT

(37)

## UNIVARIATE

NUMERIC  
CATEGORY

- COUNT - (VALUE COUNT)
- PROBABILITY

CATEGORY  
NAME

- HISTOGRAM
- KDE
- BOX PLOT

## BIVARIATE

NUMERICAL, NUMERICAL

- SCATTER
- PAIR
- HEXBIN

CATEGORY, NUMERICAL

- STRIP PLOT
- SWARM PLOT
- BOX PLOT
- BAR PLOT  
(MEAN &  
MEDIAN)