

When there is data everywhere, then we need something to visualize that data and there comes the need of this topic called Data Visualization.

To clean the data we follow a process called Dimensionality Reduction. We have certain algorithms to do this task:-

1. PCA
2. T-SNE
3. LDA

} They help us to convert large amount of data etc. into a 2D or 3D data in order to visualize them

This part of data visualization is a predominant one, in initial EDA

Exploratory Data Analysis (EDA)

process of getting data

→ There are various packages in python for data visualization, but the most prominent is Matplotlib.

→ we also use Seaborn, which is also built on top of matplotlib.

* Matplotlib Architecture:

There are '3' diff layers in this:-

- ① Backend layer
- ② Artist layer
- ③ Scripting layer

Backend layer

what is Backend layer?

This is the bottom most layer of a figure which contains implementation of several functions that are required for plotting.

3 classes of Backend layer:

→ FigureCanvas

the layer / surface on which the figure will be drawn

→ Renderer

the classification that takes care of the drawing on the surface.

→ **Event**

to handle the keyboard and mouse events

Artist layer:

note: Compared to **Artist layer**, we don't work much with the **Backend layer**.

→ this is the second layer /middle

e.g. paper is the **Figure Canvas** and
→ sketch pen as **renderer**

→ Then the hand of the painter is the **Artist layer**.

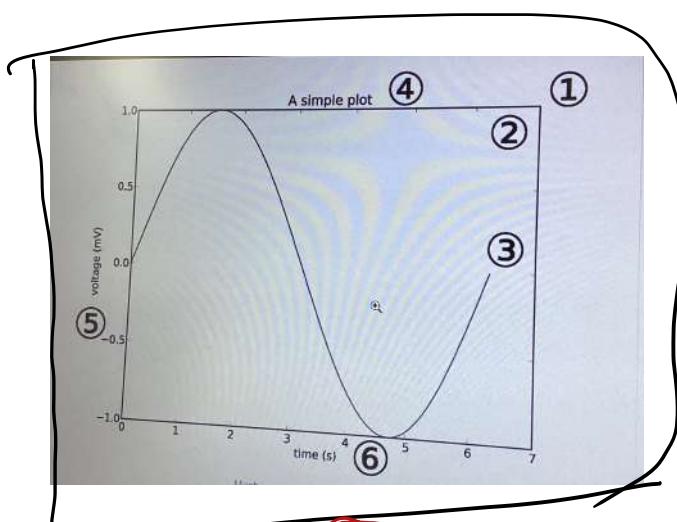
which has certain **functions**, it knows how to sketch to get the exact figure.

→ Just like in Backend layer we have 3 classes. But in Artist layer we have few classes like:-

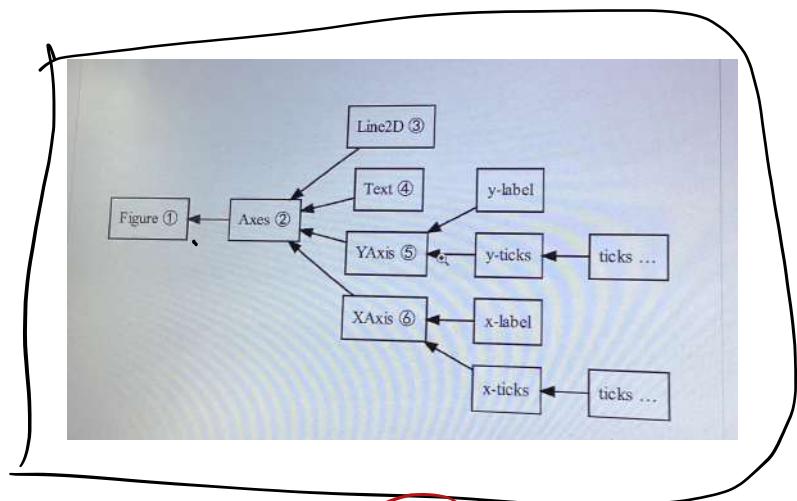
- ① Figure
- ② Axes
- ③ Axis

→ Figure is the **topmost** one

→ And we can add multiple no. of axes to a figure, upon which the plot is done.



①



②

Scripting Layer:

→ This is the top most layer on which majority of our codes will play around.

- For day to day works we almost rely on this scripting layer for ~~matplots~~.
- The methods in this layer, almost automatically takes care of layers and all we need to care about is the current state (figure & subplot).

Hence, AKA Stateful interface.

Common Terms:-

```
| import matplotlib.pyplot as plt  
| import Seaborn as sns
```

Axes: (Class that has attributes)

→ it is the entire area of a single plot in the figure.

→ it contains:

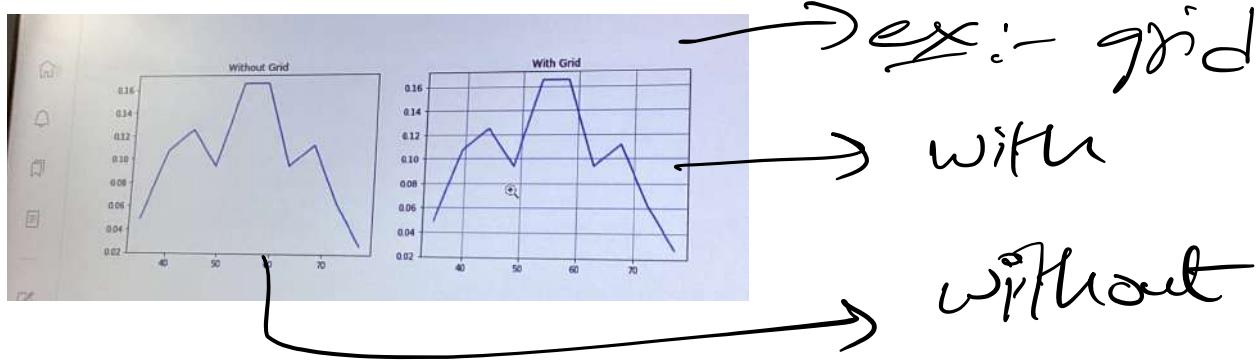
- ① adding title
- ② giving labels

③ selecting 6n values for diff types of plots. on each axes. etc..

→ we can have multiple axes in a single plot, by which we can combine multiple plots into a single figure

Grid:

- when a grid is enabled in the plt we can see horizontal and vertical lines at the background.
- This can be enabled `plt.grid()`
- used for having a rough estimation.

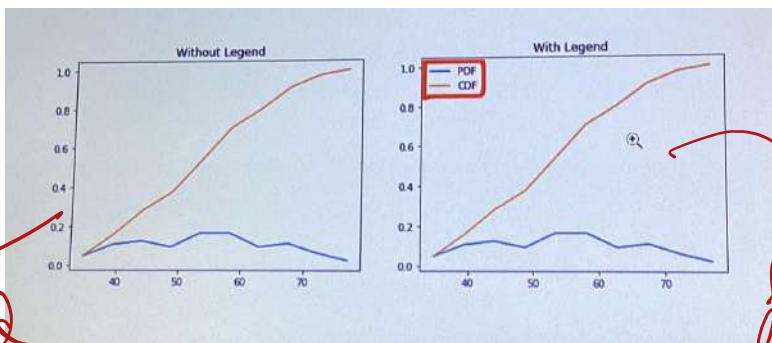


Legend:

- it is a labeled list of representation of diff diff plots which are present in a figure (on that

(particular axis)

- mostly helpful when we have diff plots in a particular figure



in this figure, legend is not there.

) in this figure, we have a legend (B)
it's helpful to view diff plots easily.

Subplots:

- when we want '2' or more plots in a single image , we can make use of subplot
- `plt.subplot(xyz)`

xyz is a '3' digit integer

X = no. of Rows

Y = no. of Columns

i = index number of that plot

e.g:-

plt. subplot(131)

1 Row

3 Columns

index value

1

note: In addition to subplot you can
also use gridspec

Title:

→ we can set a title to a plot using

plt.title()

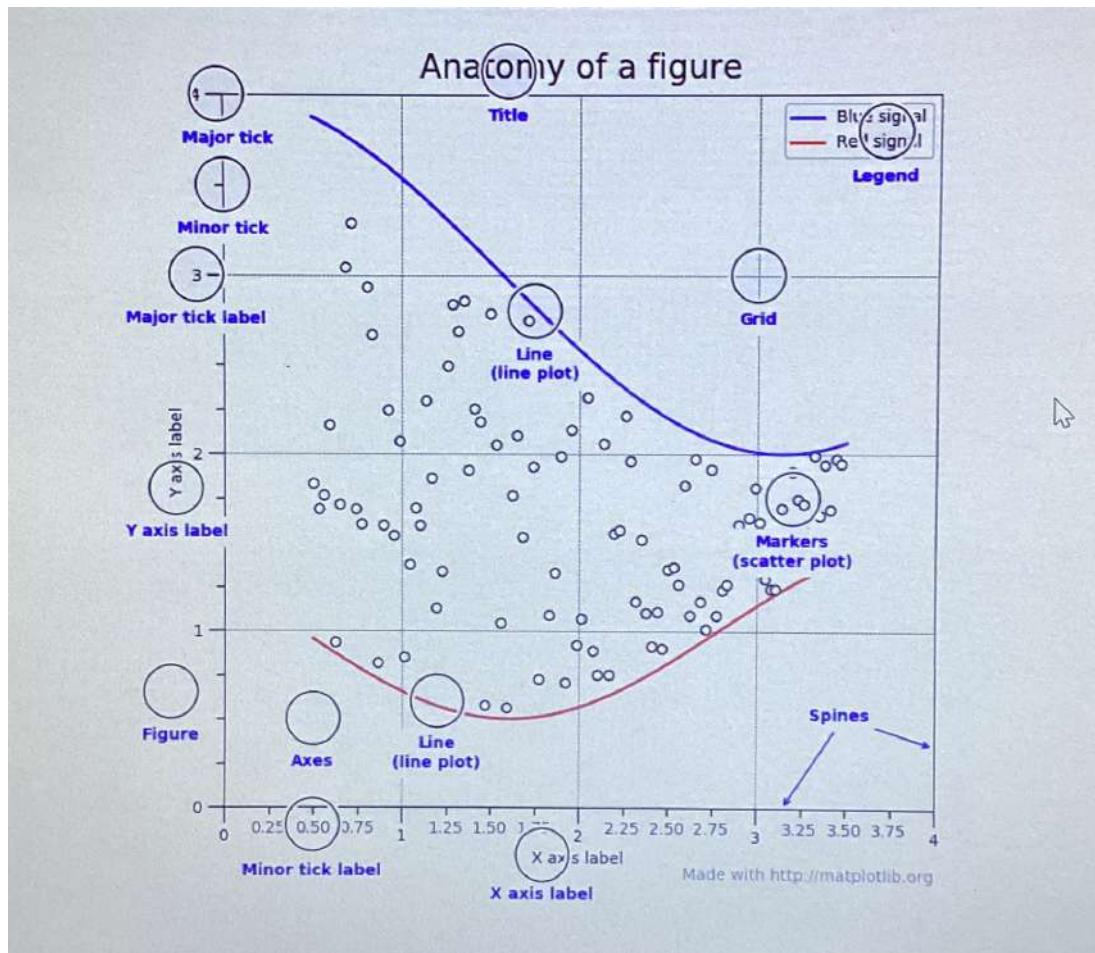
→ Xlabel: to set the label for X axis

plt.xlabel()

→ Ylabel: to set the label for Y axis.

plt.ylabel()

For a broad view, you can refer the below



Different types of analysis:-

There are '3' types of analysis

(1) Univariate: In univariate analysis we will be using a single feature to analyze almost of its properties.

(2) Bivariate: Comparing exactly b/w '2' features

③ Multivariate: Comparing more than 2' Variables

We will cover the below plots:

- ① Scatter plot (B)
- ② pair plot (M)
- ③ Box plot (U)
- ④ Violin plot (U)
- ⑤ Distribution plot (U)
- ⑥ Join plot (U & B)
- ⑦ Bar chart (B)
- ⑧ Line plot (B)

Scatter Plot :-

- one of the most commonly used plot for simple data visualization.
- available in 2D and 3D also
- 2D here is used
 - find patterns
 - clusters
 - separability of the data.

1stx

plt.scatter(x, y)

e.g. on the iris dataset we perform the scatter plot()

```
plt.scatter(iris["sepal_length"], iris["sepal_width"])
plt.xlabel("sepal length")
plt.ylabel("sepal width")
plt.title("Scatter plot on Iris dataset")
```

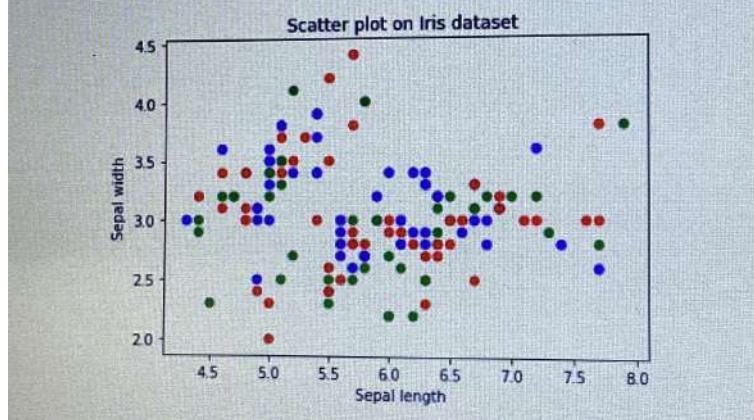
This gives the o/p having plotted the respective X and Y values.

Let's try to change the colour of the plt.

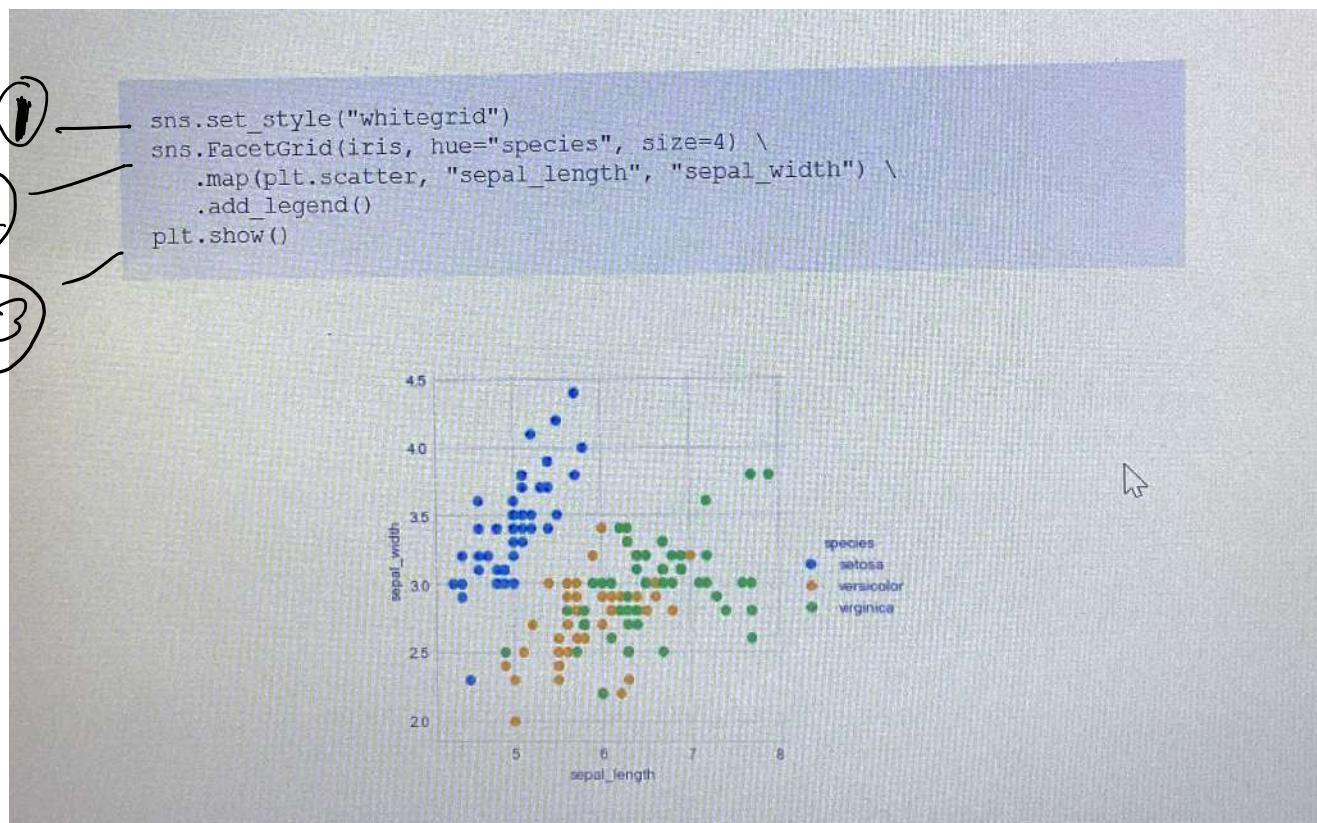
```
plt.scatter(iris["sepal_length"], iris["sepal_width"],
            color = ["r", "b", "g"])
```

```
plt.xlabel("sepal length")
plt.ylabel("sepal width")
plt.title("Scatter plot on Iris dataset")
```

→ colors are assigned to points on basis of how they were present in dataset.



- if we want colors acc to the class
- in this case, we have '3' class label
- This is where **Seaborn** pops up , with its visualisation stuff.



- ① Using seaborn , we are setting the grid upon which we plot at a "white grid"

- ② The parameter hue in the facetgrid decides the colour of each data points.
- ③ we used "Species" column in hue, so that we got this plot colored in this manner.

Note: This makes Seaborn a bit more superior than matplotlib when it comes to visualization.

⇒ for 3d scatter plots we can use plotly

Pair Plot :

→ we can scatter plot for 2D, 3D but for 4D we have to use pair plot

→ if we have ' n ' features, pair plot will create us $(n \times n)$ figure

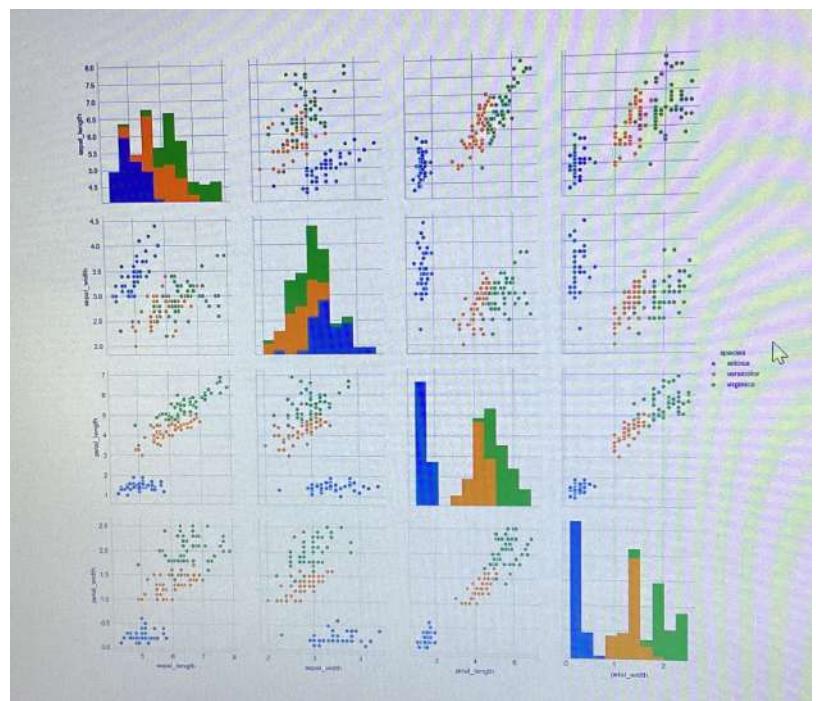
diagonal plots will be histogram plots of the feature corresponding to that row

Rest of plots are the:-

- ① Combination of feature from

each row in y-axis
feature from each column in x-axis

1
sns.set_style("whitegrid");
sns.pairplot(iris, hue="species", size=3);
plt.show()
2
3



- By performing pairplot we can observe which two features can well explain the data
- and then we can use scatter plot b/w those 2 features to explore further.

Note: Since, we will get $n \times n$ plots for n features, pairplot may become complex when we have more no. of feature.

⇒ In such case we use Dimensionality Reduction (DR) to map data into 2D plane and visualize it using our

Very own 2d Scatter plot

Box plots:-

- It can used mostly to obtain more of the **statistical details** about the data.
- The st-lines at the maximum and minimum are also called as **Whiskers**.
- Points outside of the whiskers are considered to be **Outliers**.
- Boxplot is available in Seaborn library.
- This plot actually gives us the representation in the form of **Quartiles**.
 - 25^{th}
 - 50^{th}
 - 75^{th}
- IBR**
- Inter Quartile Range
- Using box plot we can see this also !!!

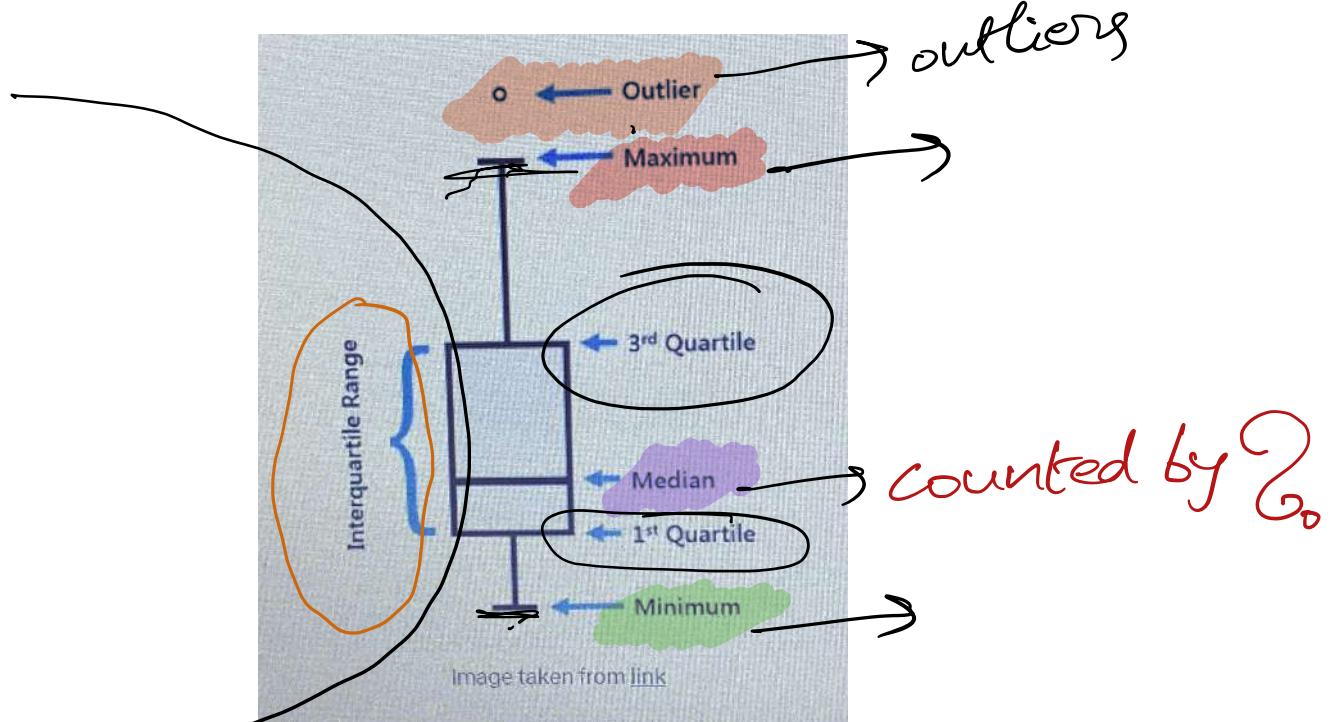
Code: $x \rightarrow$ Variable to be predicted

$y \rightarrow$ independent feature

Note: because this uses only 1 indep Variable, we call it a Univariate

⇒ Using the box plot, we can see how much of data is present in 1st Quartile and how much as Outliers.

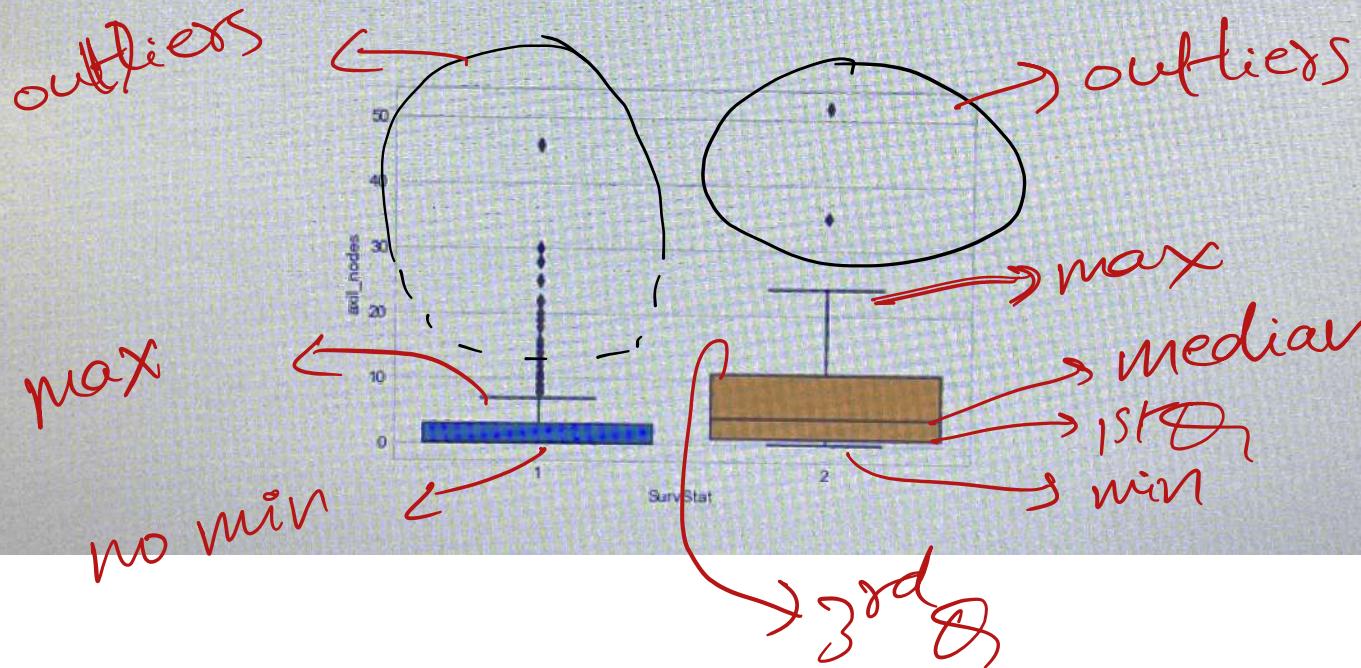
Note: We use these plots to find out outliers and to clean them before sending them to a training model. As they impact the ML process.



→ fetch → on X axis → on Y axis → the data frame we use

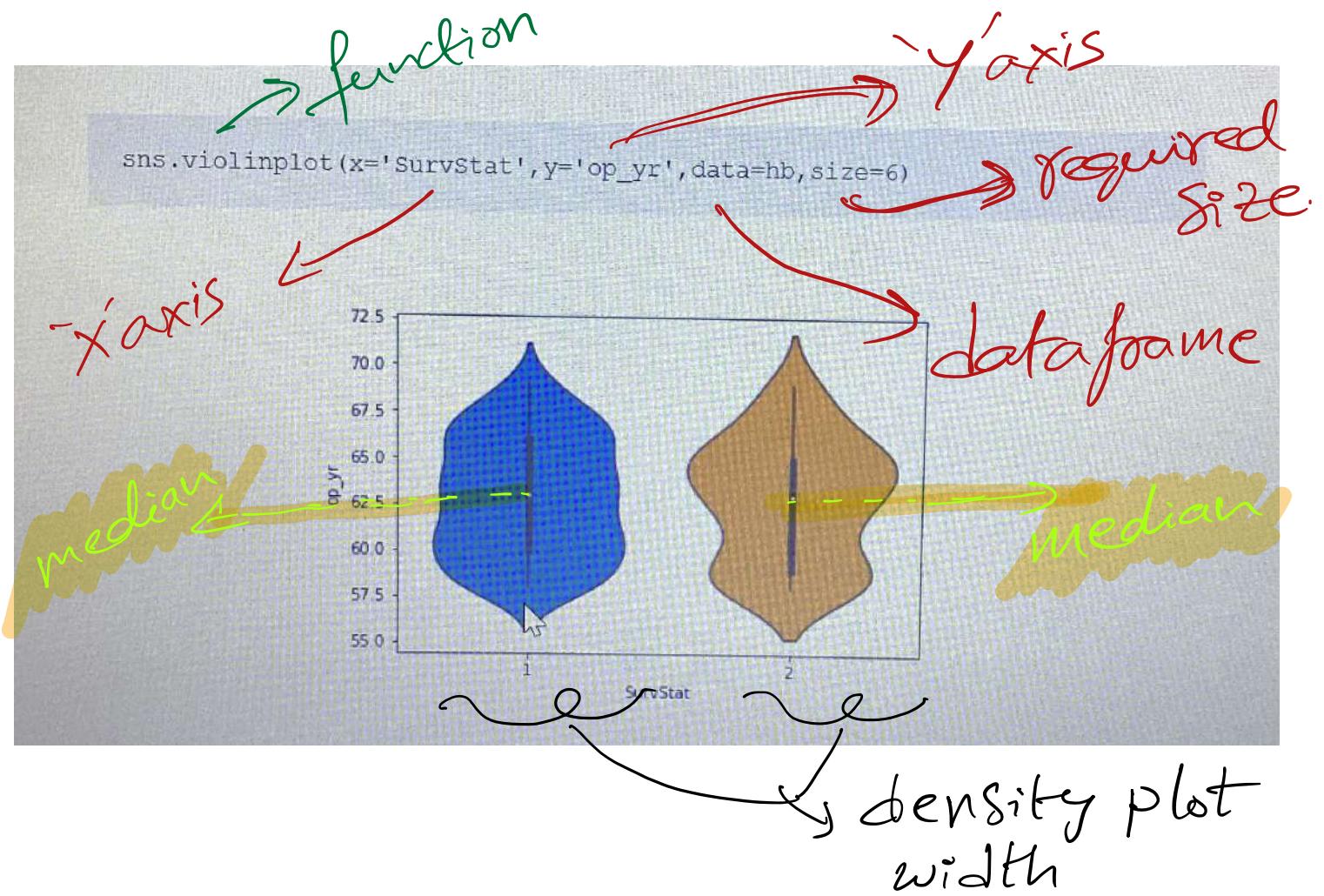
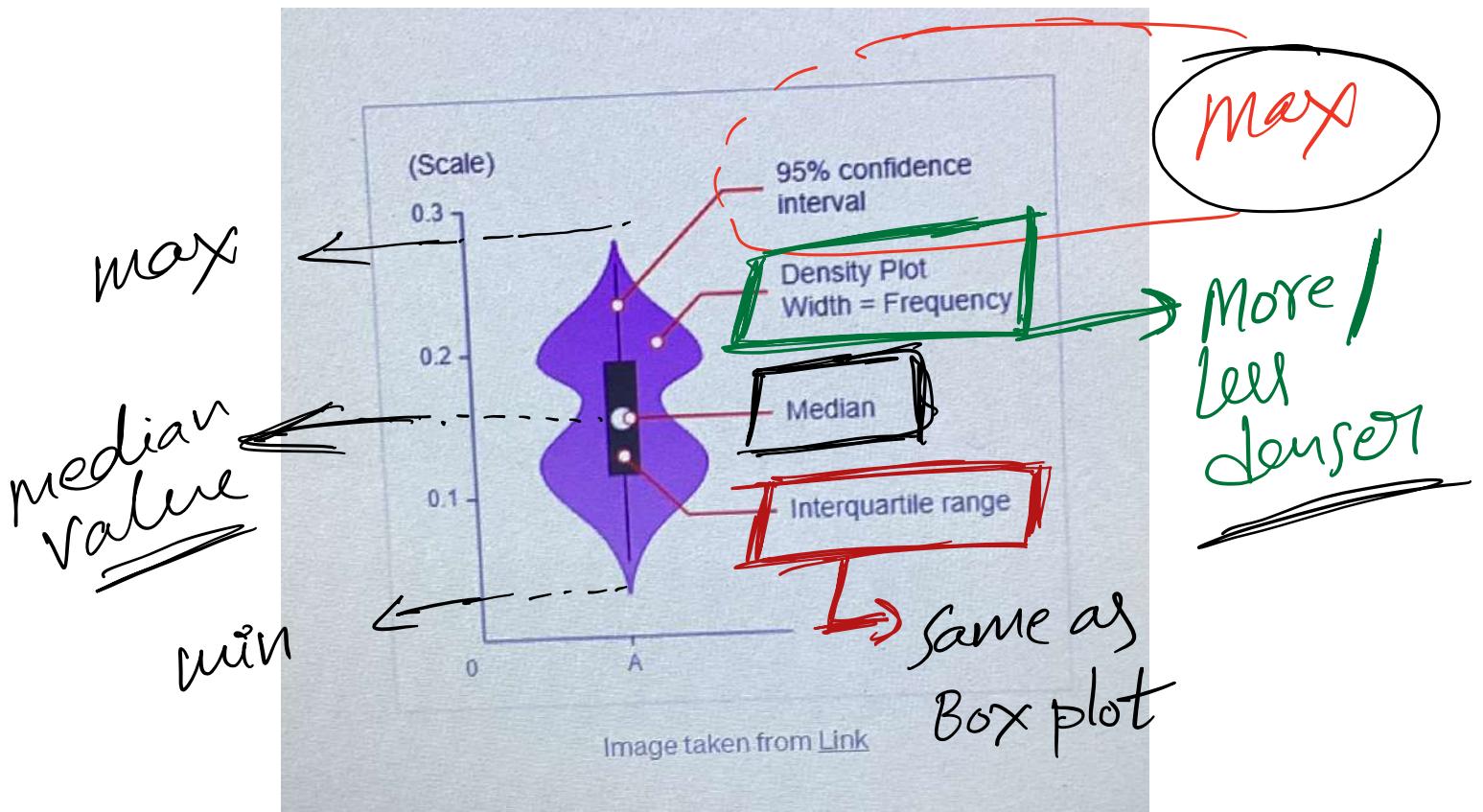
```
sns.boxplot(x='SurvStat', y='axil_nodes', data=hb)
```

Q1



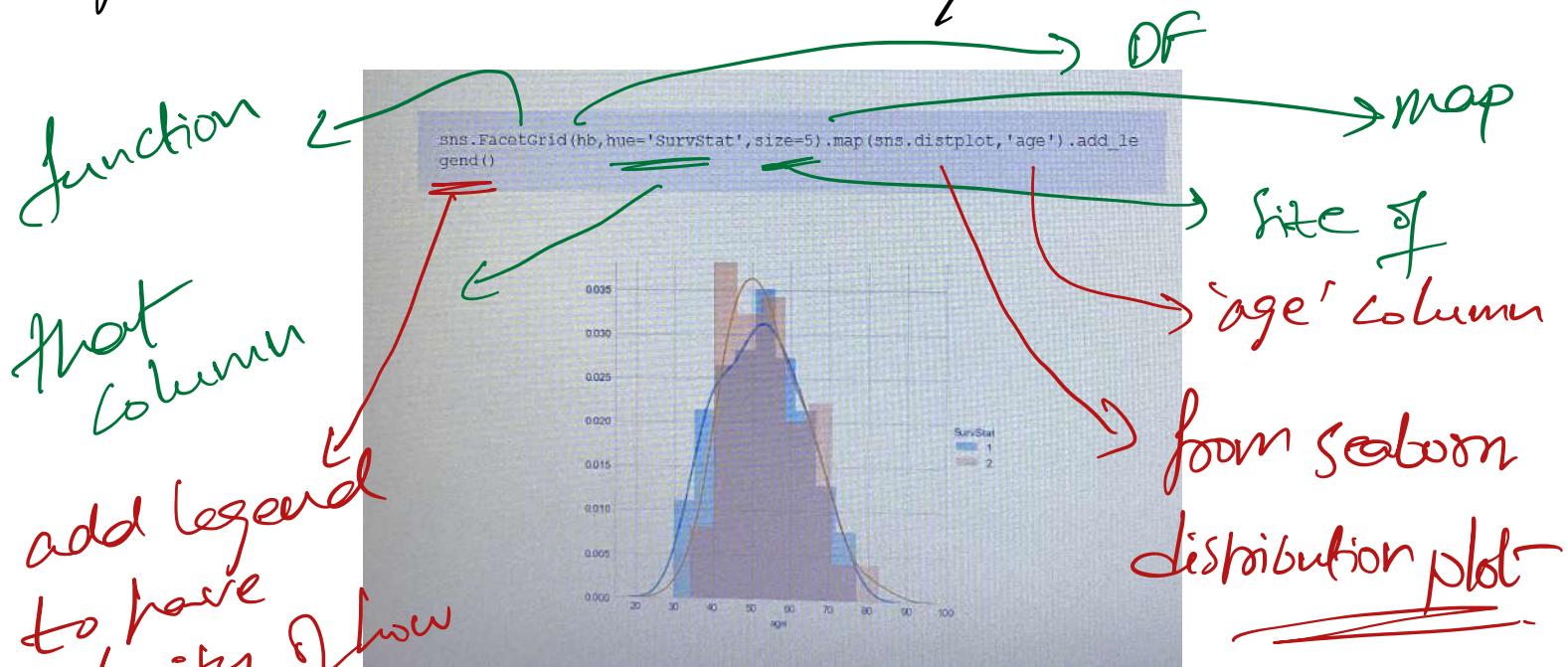
Violin plots :-

- The violin plots can be inferred as a combination of Box plot at the middle and distribution plots [Kernel Density Estimation]
- This can give us the details of distribution like whether the distribution is multimodal, skewed etc..
- gives us 95% Confidence interval.
- This plot is also from Seaborn package.



Distribution Plot:-

- one of the best univariate plot
- we use this plot when we have only one variable (most of the times)
- This plot gives us a combination of PDF + Histogram in a single feature
- The sharp block like structures are called histograms and the smooth curve is called Probability Distribution Function (PDF)
- The pdf of a curve can help us to identify the underlying distribution of that feature which is one major take away from data visualization / EDA



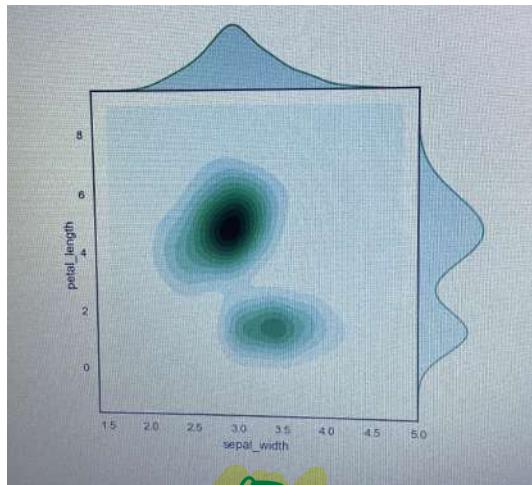
*Many
many graphs*

Join Plot:

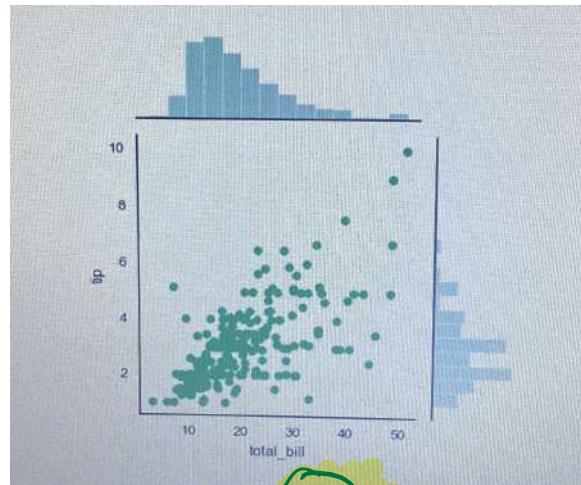
- In a single figure we can do both **single variate** and **bivariate** analysis.
- The **main** plot will give us a **bivariate** analysis,
- Whereas on the **top** and **right side** we will get **univariate** plots of both the variables that were considered.
- There is **kind** parameter in Seaborn's **join plot** function.
- The plot in the figure below is **KDE** (**Kernel Distribution Estimation**)
- More points, it will be dark and if less points it will be in light color

Note: The '2' most important plots we use will be scatter plot for bivariate and

distribution plot for univariate and since we are getting both in a single plot as shown below, it will make our work more easier.



①

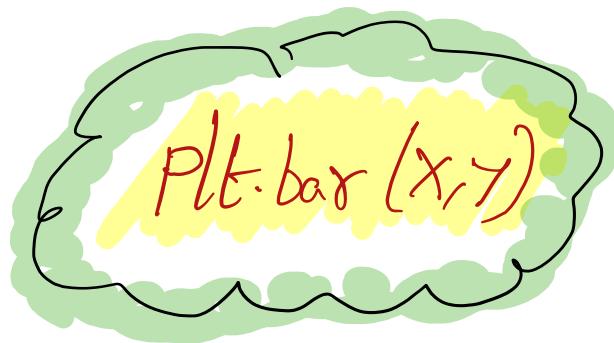


②

Bar chart:-

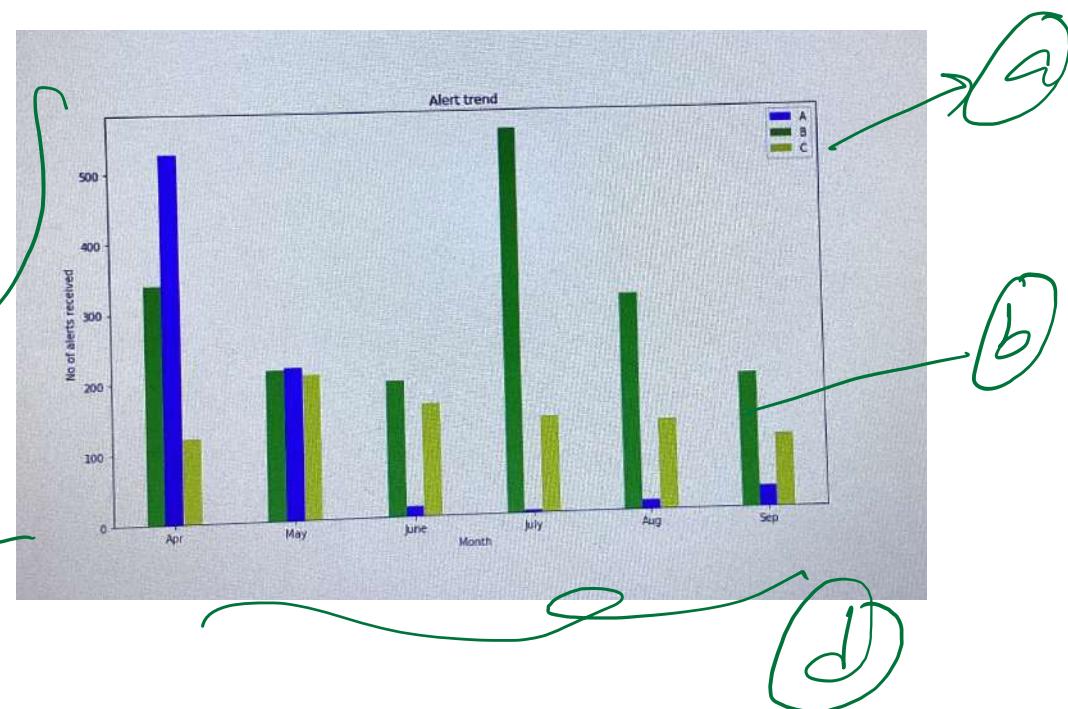
→ used to analyse data like sales figure every week, revenue from a product, no. of visitors to a site on each day of a week etc - - -

→ $X =$ Position of bars
 $y =$ length/height



→ ex:-

```
①  
② a=np.arange(6)  
w=0.15  
③  
④ fig,ax=plt.subplots(figsize=(12,7),edgecolor='k')  
p1=ax.bar(a,d,w,color='b')  
p2=ax.bar(a-w,c,w,color='g')  
p3=ax.bar(a+w,e,w,color='y')  
ax.set_xticks(a)  
⑤ ax.set_xticklabels(['Apr','May','June','July','Aug','Sep'])  
ax.set_title('Alert trend')  
ax.legend((p1[0],p2[0],p3[0]),('A','B','C'))  
plt.xlabel('Month')  
plt.ylabel('No of alerts received')  
⑥  
⑦ #plt.grid()  
⑧ plt.show()  
⑨
```

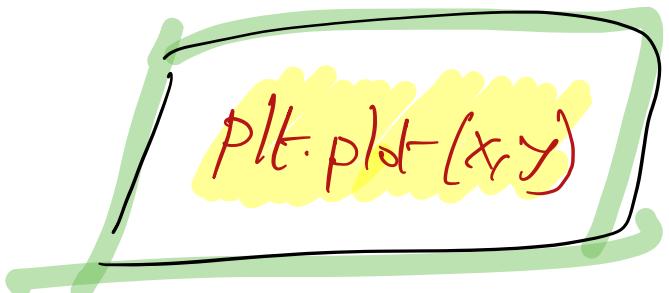


- ① Using numpy, we generate '6' nums
- ② we declare width(w) of some value
- ③ declaring subplots
- ④ for every dataset assigning proper colour
- ⑤ assigning x-axis thing like:-

- ① X ticks
 - ② X tick labels
 - ③ title
- ④ giving a legend
 - ⑤ declaring what will be on X and Y labels
 - ⑥ grid for figure
 - ⑦ show the plot

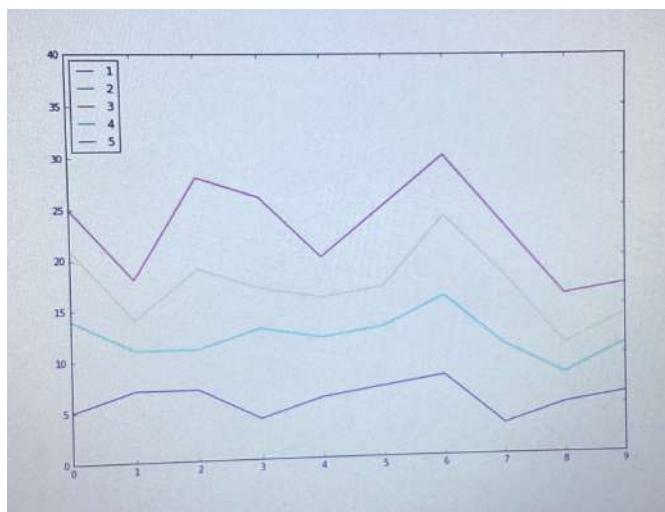
Line plots:

- You can see this plot in every nook and corner of any sort of analysis b/w '2' variables.
- graph will be connected by st. lines
- we call plot multiple lines inside a single frame as shown below where you need to add multiple plt.plot() commands with each line representing a diff color parameter
- These are used right from performing



distribution comparison using
DDplots tell CV tuning elbow method

and analysing performance of a model
using AUC curve



Note)

so far we saw some widely employed

methods that are used to extract

useful information / insight from the data.

we can also use some visualization tools
to convey our find info / inference in
form of plots to the audience.

HeatMap: (to visualize correlation pattern)

→ The heatmaps from Seaborn library will create

a grid like plot along with an optional color bar.

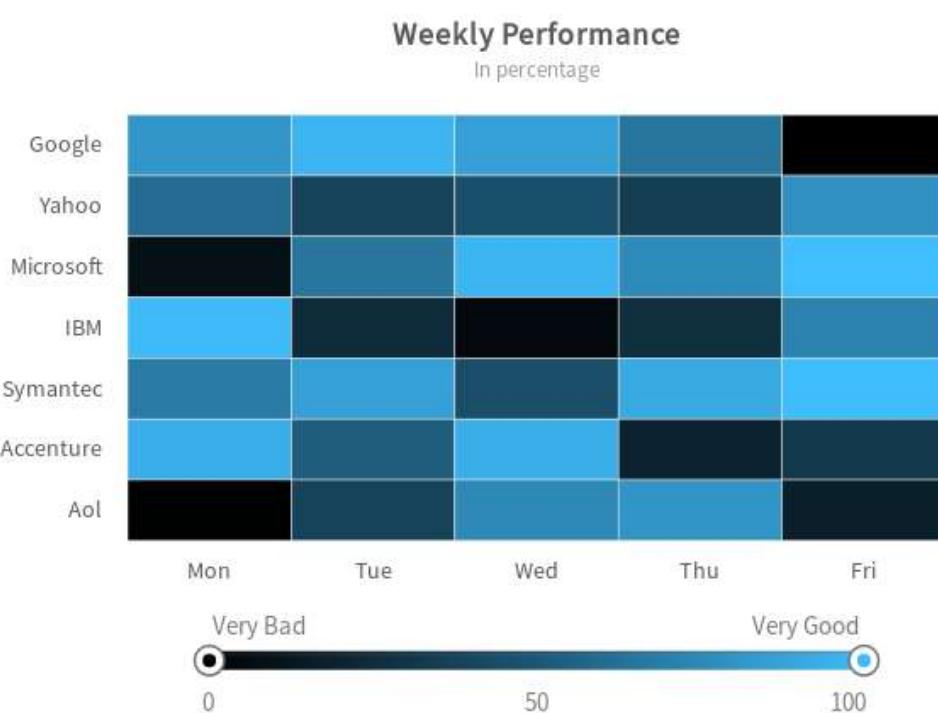
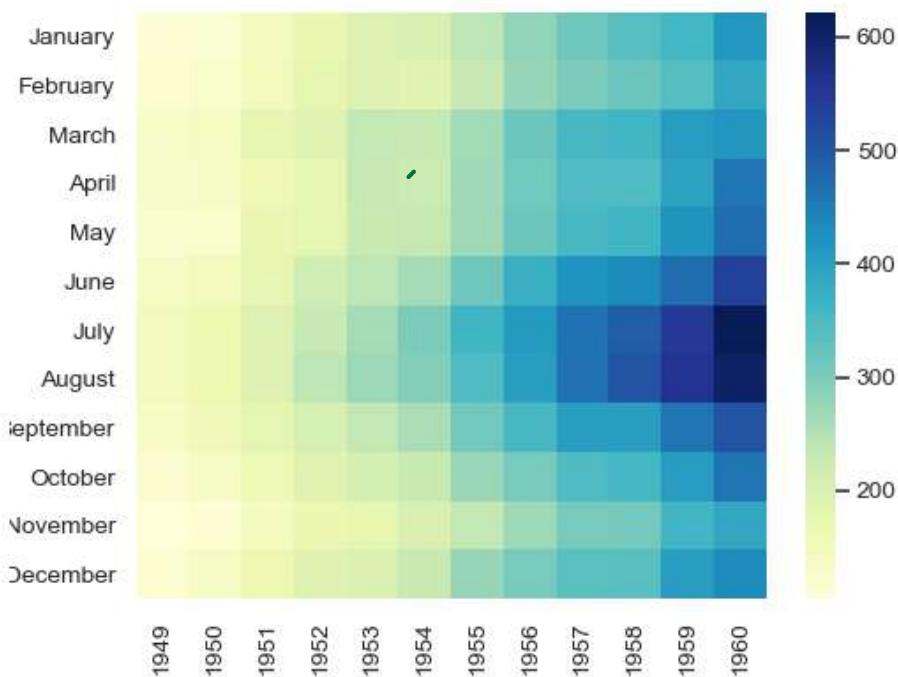
- We provide a 2D input matrix with certain values on each element to the heatmap and it exactly reproduces the output plot in the same shape as that of input matrix
- And each tile are colored based on the values provided in each elements of matrix to their corresponding tile.
- This can be much useful in cases where there will be a bigger matrix and we want to find which value is higher, lower and so on.... by simply looking at the diff colors tones used.

e.g. weekly performance chart of different companies can be plotted using heatmaps.

Note:

In ML Applications, it can be used in representing confusion matrix of a model, used in hyperparameter tuning to plot

error values b/w '2' diff hyperparameters
ek ---



Wordcloud:

→ It is nothing but creating an image

that contains all the words in a passage/string with diff size, as per the frequency of occurrence in that passage.

→ The word will appear bigger and bolder if it occurs more no. of times in a given passage.

Note 

- ⇒ we can conclude like what's discussed in a paragraph without reading it much.
- ⇒ The most repeating word will appear bold and bigger than the least occurring one.



It can be used a lot with text data analysis like, finding the common topics in a cluster while performing sentiment analysis and SEO optimizations etc.--

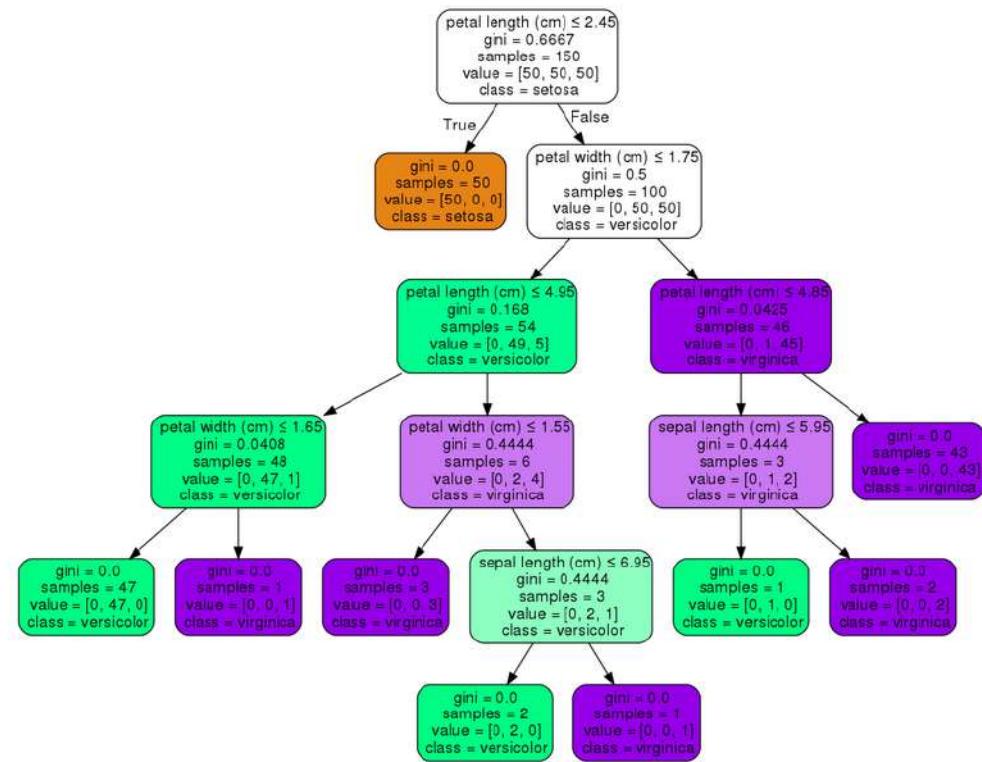
below is the representation of
"Word cloud"



Graphviz:

- The Decision tree algorithms are one of the popular non-linear model.
 - If you want to view a linear model like linear regression you can do it simply using matplotlib / seaborn
 - whereas to visualize trees, we use a special tool called **Graphviz**
 - The sklearn has a good implementation of a function called **export graphviz** which can help us to convert the decision tree model into dot format,

which is supported by graphviz tool.



Dendograms:

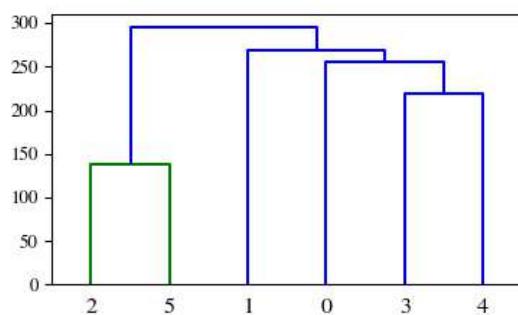
→ This is the method used to visualize the clusters formed when using the hierarchical clustering technique like :

Agglomerative clustering

→ The below figure is an example for dendrogram. Here each chain/link b/w

points, means they fall under same cluster.

Note: The dendrogram structure may be too much complex when the number of datapoints are large.



Histogram:-

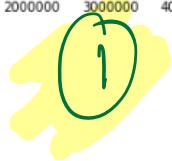
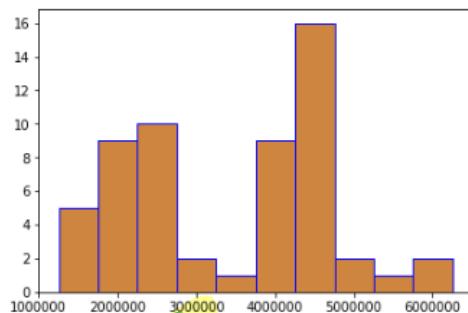
- A histogram takes in a series of data and divides the data into a number of bins.
- It then plots the frequency data points in each bin (i.e. interval of points)
- It is useful in understanding the count of data ranges.
- When to use:

① We should use histogram when we need the count of the variable in a plot.

→ The Cumulative property gives us the end added value and helps us understand the increase in value at each bin.

Histogram

```
In [3]: 1 plt.hist(np_data['GrandCanyon'],
2             facecolor='peru',
3             edgecolor='blue',
4             bins=10)
5 plt.show()
```

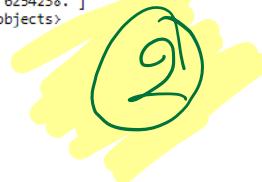


Components of a histogram

- n: Contains the frequency of each bin
- bins: Represents the middle value of each bin
- patches: The Patch object for the rectangle shape representing each bar

```
In [4]: 1 n, bins, patches = plt.hist(np_data['GrandCanyon'],
2                               facecolor='peru',
3                               edgecolor='blue',
4                               bins=10)
```

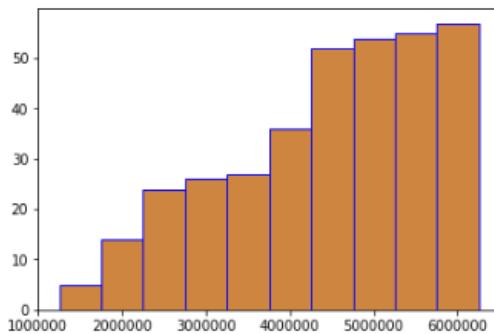
```
n: [ 5.  9. 10.  2.  1.  9. 16.  2.  1.  2.]
bins: [1253000. 1753123.8 2253247.6 2753371.4 3253495.2 3753619. 4253742.8
4753866.6 5253990.4 5754114.2 6254238.]
patches: <a list of 10 Patch objects>
```



The cumulative property

If True, then a histogram is computed where each bin gives the counts in that bin plus all bins for smaller values. The last bin gives the total number of datapoints.

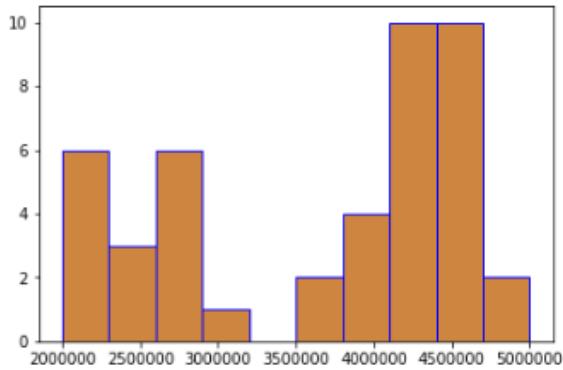
```
In [5]: 1 plt.hist(np_data['GrandCanyon'],
2             facecolor='peru',
3             edgecolor='blue',
4             bins=10,
5             cumulative=True)
6 plt.show()
```



Check the histogram to a range of values

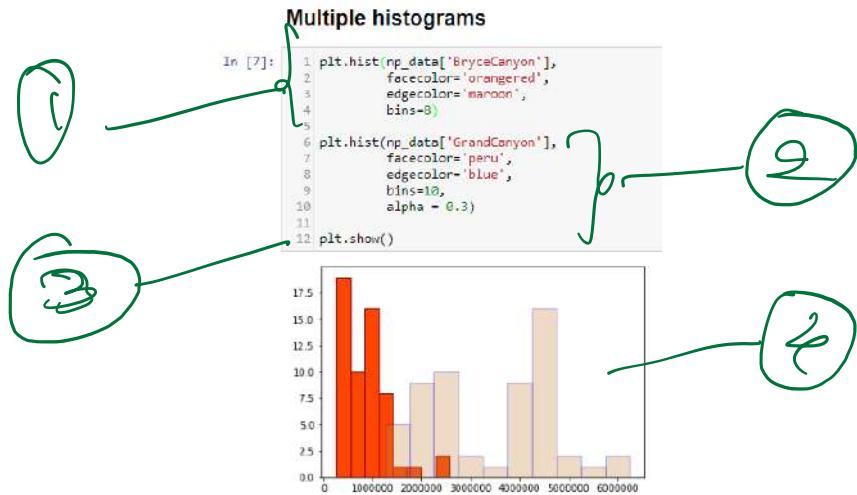
We only look at the data points within the range 2M-5M. This realigns the bins in the histogram

```
In [6]: 1 plt.hist(np_data['GrandCanyon'],
2           facecolor='peru',
3           edgecolor='blue',
4           bins=10,
5           range=(2000000, 5000000))
6
7 plt.show()
```



④

- Multiple histograms are useful in understanding the distribution between 2 entity variables.
- we can see that Grand Canyon has comparably more visitors than Bryce Canyon.



Pie chart:

- It is a circular plot which is divided into slices to illustrate numerical proportion.
- The slice of a pie chart is to show the proportion of parts out of a whole.
- When to use: pie chart should be seldom used as it is diff to compare sections of the chart.

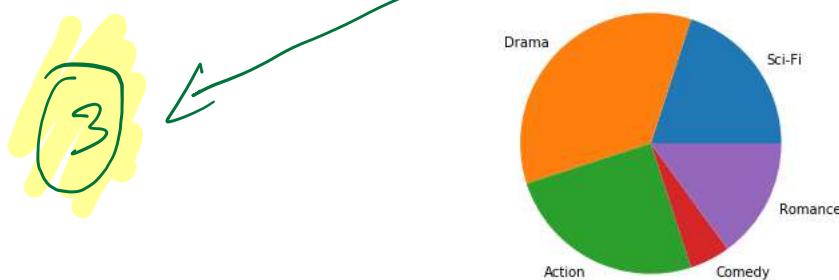
② Bar plot is used instead as comparing sections is easy.

Pie Chart

Pie chart shows distribution of data based on proportion of pie occupied.

In [3]:

```
1 plt.pie(t_mov['Percentage'],
2          labels=t_mov['Sector'])
3
4 plt.axis('equal')
5
6 plt.show()
```



Pie chart components

- wedges: A list of Patch objects representing each wedge
- texts: List of Text objects representing the labels
- autotexts: List of Text objects for the numeric values - this is only available if the autopct value for the pie chart is not None

```
In [4]:  
1 wedges, texts, autotexts = plt.pie(t_mov['Percentage'],  
2 labels=t_mov['Sector'],  
3 autopct='%.2f')  
4  
5 plt.axis('equal')  
6  
7 print('Wedges:', wedges)  
8 print('Texts:', texts)  
9 print('Autotexts:', autotexts)  
10  
Wedges: [, <matplotlib.patches.Wedge object at 0x000001F480051B00>, <matplotlib.patches.Wedge object at 0x000001F4800582E  
he>, Wedge object at 0x000001F480059A00>]  
Texts: [Text(0.889919, 0.546564, 'Sci-Fi'), Text(-0.777817, 0.777817, 'Drama'), Text(-0.49939, -0.980107, 'Actio  
n'), Text(0.988107, 'Comedy'), Text(0.488541, 0.352671, 'Romance')]  
Autotexts: [Text(0.488541, 0.352671, '28.00'), Text(-0.424264, 0.424264, '35.00'), Text(-0.272394, -0.534684, '25  
4, -0.534684, '5.00'), Text(0.534684, -0.272394, '15.00')]
```

Pie chart customizations

We set the values for the colors and autopct properties. The latter sets the format for the values to be displayed.

```
In [5]:  
1 colors = ['darkorange', 'sandybrown', 'darksalmon', 'orangered', 'chocolate']  
2  
3 plt.pie(t_mov['Percentage'],  
4 labels=t_mov['Sector'],  
5 colors=colors,  
6 autopct='%.2f')  
7  
8 plt.axis('equal')  
9  
10 plt.show()
```

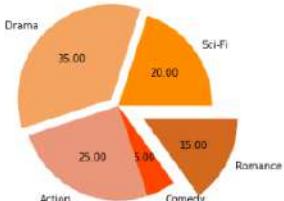


The explode property

To highlight a particular wedge of the pie chart, we use explode to separate it from the rest of the chart.

The value for "explode" represents the fraction of the radius with which to offset each wedge.

```
In [6]:  
1 explode = (0, 0.1, 0, 0, 0.3)  
2  
3 plt.pie(t_mov['Percentage'],  
4 labels=t_mov['Sector'],  
5 colors=colors,  
6 autopct='%.2f',  
7 explode=explode)  
8  
9 plt.axis('equal')  
10  
11 plt.show()
```



note:

There many different types of plots, acc to the requirements of the project we prefer 1st plot to use.

