

## Module 2

### Matplotlib

- Visualizing data is crucial to quickly understanding trends and relationships in your dataset.
- Matplotlib is one of the most popular libraries for plotting with Python.
  - Sometimes known as "Grandfather of plotting" and visualization libraries.
  - Many visualization libraries are directly built on top of Matplotlib (e.g. Seaborn and Pandas built-in visualization).
  - Matplotlib is heavily inspired by the plotting function of MATLAB programming language -
  - It also allows creation of almost any plot type and heavy customization
    - Two separate approaches to creating plots -
      1. Functional plots.
      2. GOP based plots.
    - Matplotlib Basic
      - Functional Methods
      - Matplotlib figures and subplots
        - GPP Method

## Topic covered:-

- matplotlib basic and functions
- Matplotlib Figures
- Matplotlib subplots
- Matplotlib styling
- Exercise Question and Solution.

Two main goals with matplotlib:

- Functional Relationship

$$y = 2x$$

- Be able to plot out a relationship between raw data points:

$$x = [1, 2, 3, 4]$$

$$y = [2, 4, 6, 8]$$

## 13 matplotlib Basic

- Basic way to use Matplotlib is functional plot calls:
  - plt.plot(x,y)
- These function calls are really simple to use, but they don't actually allow high degree of use.
- need plt.show() in .py scripts but in jupyter notebook not needed.

### Code #

```
import matplotlib.pyplot as plt  
%matplotlib inline
```

```
import numpy as np  
x = np.arange(0, 10)  
y = 2**x
```

```
plt.plot(x, y)  
plt.show() — necessary in .py
```

→ ~~print~~ doesn't show the address on the ~~print~~ memory.

[<matplotlib.lines.Line2D at 0x7066a1 ]

```
plt.plot(x, y)  
plt.title('Spring Title')  
plt.xlabel('X')  
plt.ylabel('Y')
```

add limits

plt.plot(x, y)

plt.xlim(0, 6) — set limits on scale X

plt.ylim(0, 15) — set limits on

plt.savefig('myfirst.png') save y

Saving the plot:- or provide the

path  
help(plt.savefig)

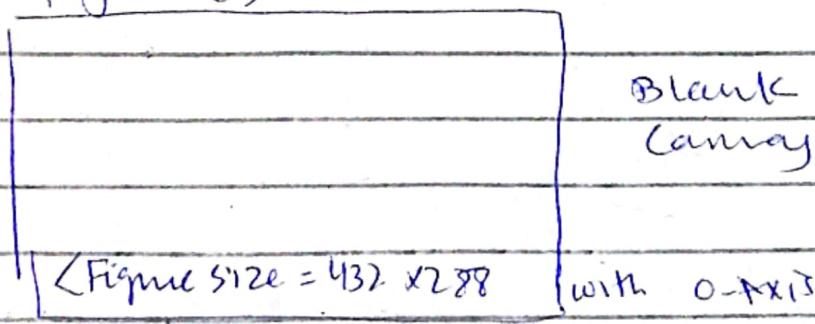
plt.savefig("C:\Users\Helly\myfirst.png")

### 03- Matplotlib - Understanding the Figure Object

part - One

Code # invisible

plt.figure()



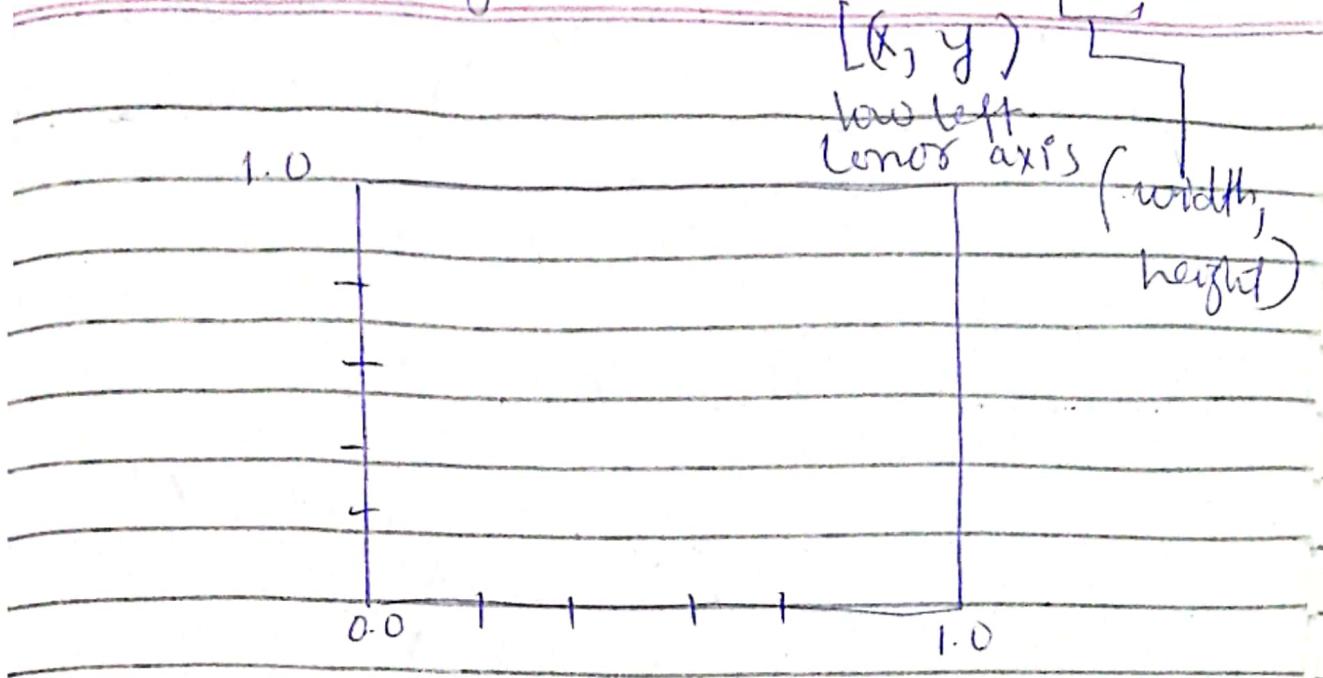
plt.figure(figsize=(10, 10))

then

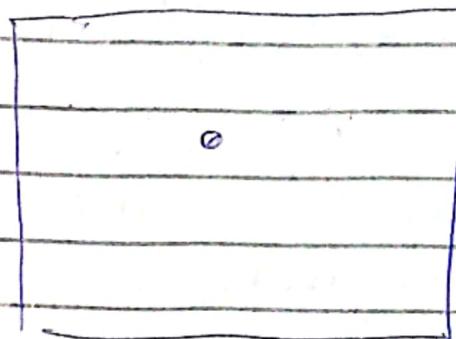
Figure size = 720 x 720

`fig = plt.figure()`

`axes = fig.add_axes([0, 0, 1, 1])`



to start

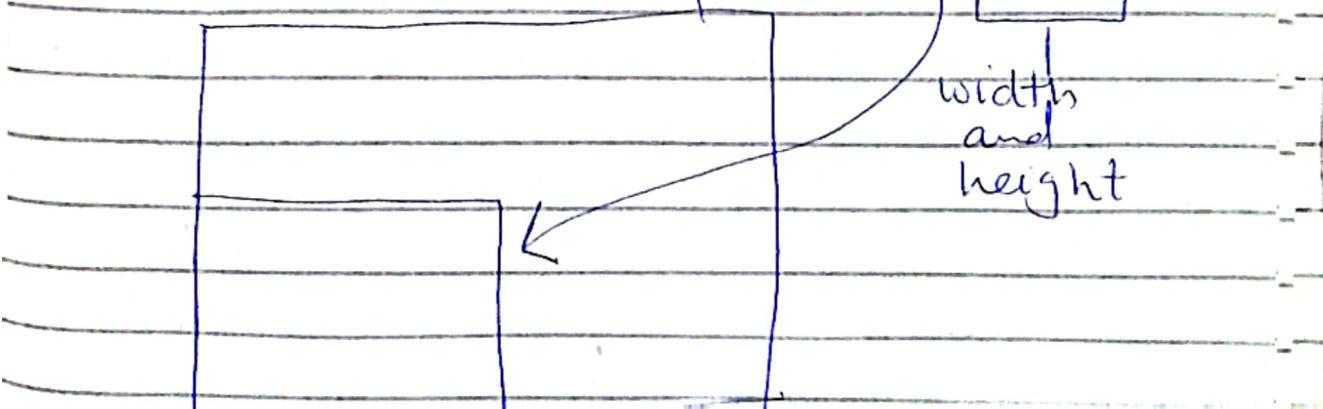


`fig.add_axes([0.5, 0.5, 1, 1])`

lower  
left  
corner of  
Axes

`fig = plt.figure()`

`axes = fig.add_axes([0, 0, 0.5, 0.5])`



axes.plot(x, y) — plot the figure

## 4 - Matplotlib - Implementing Figure and Axis part - two

a = np.linspace(0, 10, 11)

b = a\*\*4

x = np.arange(0, 10)

y = a \* 2\*\*x

fig = plt.figure()

axes = fig.add\_axes([0, 0, 1, 1])

axes.plot(x, y); axes.plot(a, b)

plt.show()

multiple axis on same figure #

fig = plt.figure()

# Large axis

axes1 = fig.add\_axes([0, 0, 1, 1])

axes1.plot(x, y)

# Small plot

axes = fig.add\_axes([0.2, 0.2, 0.5,

0.5]))

```
axes1.set_xlim(0, 50)
axes1.set_ylabel('B')
axes1.set_xlabel('Label')
axes2.plot(x, y)
```

plt.show()

copy paste  
for axes2

we can plot zoom in in the figure

## 5. Matplotlib Figure Object part 3 - Figure parameters-

Code #

```
fig = plt.figure()
axes1 = fig.add_axes([0, 0, 1, 1])
```

clear the image  
high resolution

```
fig = plt.figure(figsize=(10, 3), dpi=200)
fig.savefig('new-figure.png')
```

↓  
change

figsize

↓  
clear  
resolution

```
fig.savefig("image.png", bbox_inches='tight')
```

## 6 - Matplotlib Subplots

- Subplot functionality

- The plt.subplot() call allows us to easily create a Figure and Axes objects in side by side formation -

- The plt.subplots() command returns a tuple containing the Figure canvas and then a numpy array holding the axes object.

`fig, axes = plt.subplots(rows=1, ncols=2,`



- plt.subplots() return a tuple which by common convention we label: (fig, axes);

- fig

- This is the entire figure canvas -

- axes

- This is a numpy array

holding each of the axes according to position in the overall canvas.

Code #

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

```
a = np.linspace(0, 10, 11)  
b = a**4
```

```
x = np.arange(0, 10)  
y = 2**x
```

```
fig, axes = plt.subplots(rows=1, cols=2)  
axes[0].plot(x, y)  
axes[1].plot(a, b)
```

```
fig, axes = plt.subplots(2, 2)  
axes[0][0].plot(x, y)  
axes[1][0].plot(x, y)
```

loop through axes #

```
fig, axes = plt.subplots(3, 1)  
for ax in axes:  
    ax.plot(x, y)
```

~~Everything too close together~~

at the end of

plt.tight\_layout()

automatically  
attempt to  
space

fig, axes = plt.subplots(2, 2)

axes[0][0].plot(x, y)

axes[0][1].plot(x, y)

axes[1][0].plot(x, y)

axes[1][1].plot(x, y)

fig.subplots\_adjust()

wspace = 0.9, hspace = 0.9

axes[1][0].set\_ylabel('Y label')

axes[1][1].set\_title('Title')

fig.suptitle("Figure text", fontsize=16)

fig.savefig('new\_subplot.png',  
bbox\_inches='tight')

fig.set\_figwidth(10)

# 7. Matplotlib Styling - Legends

## Part one: Legends

Code #

```
fig = plt.figure()
ax = fig.add_axes([0, 0, 1, 1])
ax.plot(x, x)
ax.plot(x, x**2)
```

```
ax.plot(x, x**2, label = "X vs X")
ax.plot(x, x**2, label = "X vs X^2")
ax.legend()
```

```
ax.legend(loc = "lower left")
ax.legend(loc = 4)
```

associate

```
ax.legend(loc = (1.1, 0.5))
```

## B- Matplotlib Styling :-

```
x = np.linspace(0, 11, 10)
```

```
fig = plt.figure()
ax = fig.add_axes([0, 0, 1, 1])
ax.plot(x, x)
ax.plot(x, y, color = 'purple')
```

`ax.plot(x, y, color = "#63EES4")`

`ax.plot(x, x+1, color = 'purple', label = "X vs X+1")`

line width #

`ax.plot(x, x+1, color = "#e6c0ff", lw = 10)`

line width

`ax.plot(x, x+1, color = " ", linewidth = 0.5)`

linestyle #

`ax.plot(x, y, linestyle = '--')`

`ax.plot(x, y, linestyle = ":" )`

Custom linestyle #

`lines = ax.plot(x, x+1, color = 'purple', linewidth = 5)`

`lines[0].set_dashes([5, 2, 10, 2])`

## Marker: $\text{eff}$

`ax.plot(x, x+1, color='purple', lw=2,  
ls='--', markers='o')`

$$\text{markersize} = 20 / ms = 20$$

`ax.plot(x, y, color='purple', lw=2, ls='--',  
markers='o', markeredgecolor='red', markerfacecolor='orange',  
markerfacecolor='orange')`

## 9. Advanced Matplotlib Command (optional)

- matplotlib is a huge library!

Advance topics separate Jupyter  
file provided by the instructor

# Matplotlib Exercise

$$E = mc^2$$

1)  $m = \text{even}$  at evenly spaced value

$$m = [0, 1, 2, 3, 4, 5, \dots]$$

2) plot  $E = mc^2$

x = mass in grams

y = Energy of joules

3) plot the cell on logarithms

Two #

$$\text{labels} = [1\text{ Mo}, 3\text{ Mo}, 6\text{ Mo}, 11\text{ Yr}, \\ 2\text{ Yr}, 3\text{ Yr}, 15\text{ Yr}, 7\text{ Yr}, \\ 10\text{ Yr}, 120\text{ Yr}]$$

$$\text{July 16 - 2007} = [4.75, 4.98, 5.08, \dots]$$

$$\text{July 16 - 2007} = [0.12, 0.11, 0.13, \\ 0.14, 0.16, 0.17, \\ 0.18, 0.46, 0.69, 1.32]$$

# Seaborn

- Seaborn is a statistical plotting library that is specifically designed to interact well with pandas Dataframe to create common statistical plot types.

- Seaborn is built directly off Matplotlib, but uses a simpler "one-line" syntax.
- Use OneLine function

```
sns.scatterplot(x='salary', y='sales',  
                  data=df)
```

## Section Topics:

- Scatter plot
- Distribution plot
- Categorical plot
- Comparison plots
- Seaborn plots
- Matrix plots

## 2- Scatterplot with Seaborn

- Scatterplot shows the relationship between two continuous features.
- Recall that continuous features are numeric variable that can take any number of values between any two values.
- Continuous features Examples:-
  - Age
  - Height
  - Salary
  - Temperature
  - prices

- Scatterplot line up a set of two continuous features and plot them out as coordinates.

Code #

```
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
import seaborn as sns  
plt.figure(figsize=(12,4))  
sns.scatterplot(x='salary', y='sales', data=df)
```

hue = 'level of education',  
either by palette = 'Dark2',  
continuous or size = 'Salary'  
categorical s = 'P200', alp1

## parameters in Scatterplot

$x = 'SI'$  display on x axis

$y = 'em'$  display on y axis

$data = df$  data which is used

$hue = 'worker'$  differentiate color based on hue

$size = 1$  size of scatter points changed based on size

$s = 100$  size of all scatter points

$\alpha = 1$  visibility of scatter points  
Transparency

$style = 'level'$  - will set markers of different shape based on category  
(can be used with)

$hue = 'level'$  for better visualization

palette = 'Dark2', colormaps in matlab docs  
check

`plt.savefig("myplot.png")`

`plt.show()` - in .py files

## 2. Distribution plots part one

### Understanding plot types

Distribution plots display a single continuous feature and help visualize properties such as deviation and average value.

There are 3 main distribution plot types:

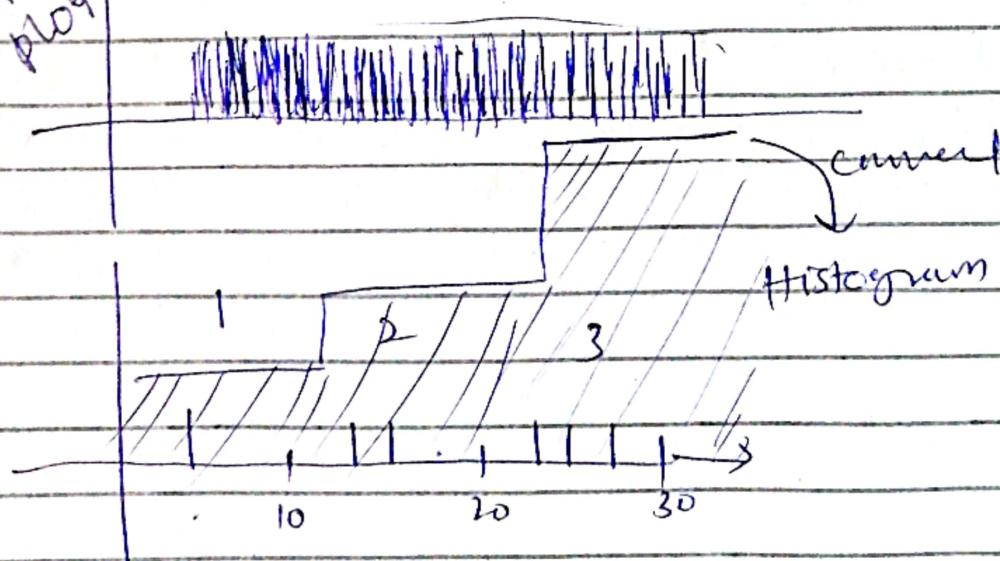
- Rug plot
- Histogram
- KDE plot

Rug plot:-

A rug plot is the simplest distribution plot and merely adds a dash or tick line for every single value.

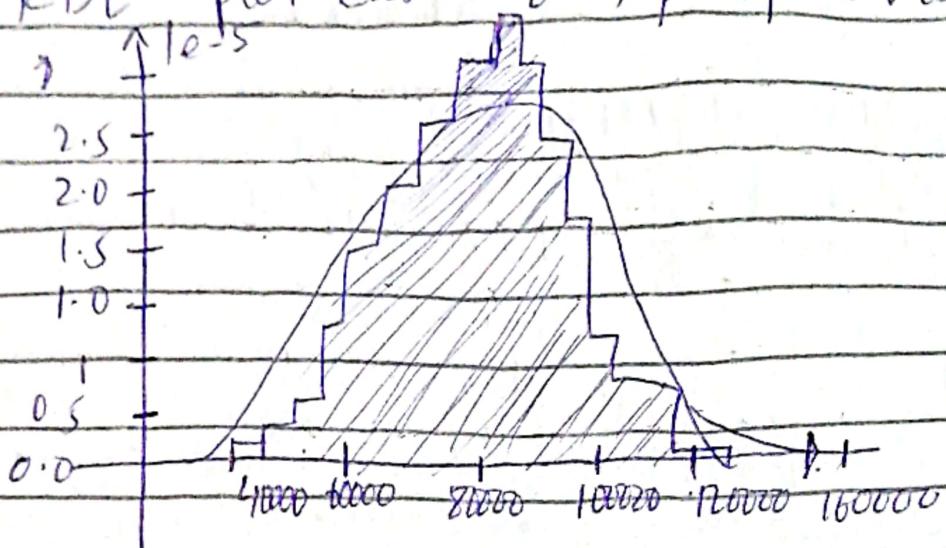
The y-axis doesn't have any meaning

rug plot



Y-axis can be seen as percent

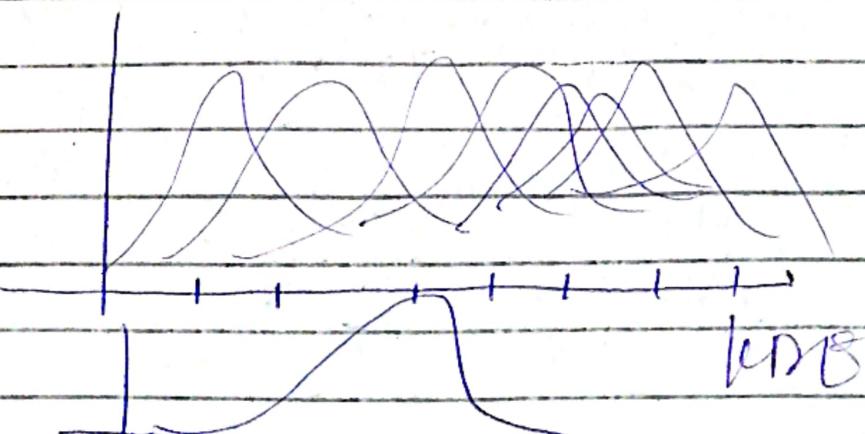
Seaborn allows us to add on a KDE plot curve on top of a histogram.



which show you what a continuous

- KDE stand for Kernel Density Estimation.
- Its a method of estimating a probability density function of a random variable.  
~~variable~~ In simple terms, its way of estimating a continuous probability curve for finite data sample.

kernel = "Gaussian distribution"



## 4- Distribution plot part - 2

Code #

#

```
sns.mugplot(x='SalaryColumn', data=df,  
             height=0.5)
```

# Do not use DISTPLOT

```
sns.set(style='whitegrid')
```

```
sns.displot(data=df, x='Salary', bins=10)
```

Multiplot color as

```
color = 'blue'
```

```
edgecolor = 'blue',
```

```
linewidth = 5,
```

```
kde = True
```

```
mug = True
```

```
sns.histplot(data=df, x='Salary')
```

only KDE plot:-

```
sns.kdeplot(data=df, x='Salary')
```

```
np.random.seed(42)
```

```
sample_ages = np.random.randint(0, 100, 20)
```

```
samples_ages = pd.DataFrame
```

```
sns.mugplot(data=df, x='age')
```

```
sns.displot(data=df, mug=True)
```

```
sns.kdeplot(data = sample, x=xage,  
            clip=[0,100], bw_adjust  
                  = 0.01)  
Shade = True
```

clip = [0, 100]

bw\_adjust = 0.01

Shade = True +

## 5. Categorical plot

statistical Estimation within categories -

- The categorical plots discussed here will display a statistical mean value per category -
- For example mean value per category or a count of the numbers of rows per category -
- Its a visualization equivalent of a groupby call() -

Two main types of plot:-

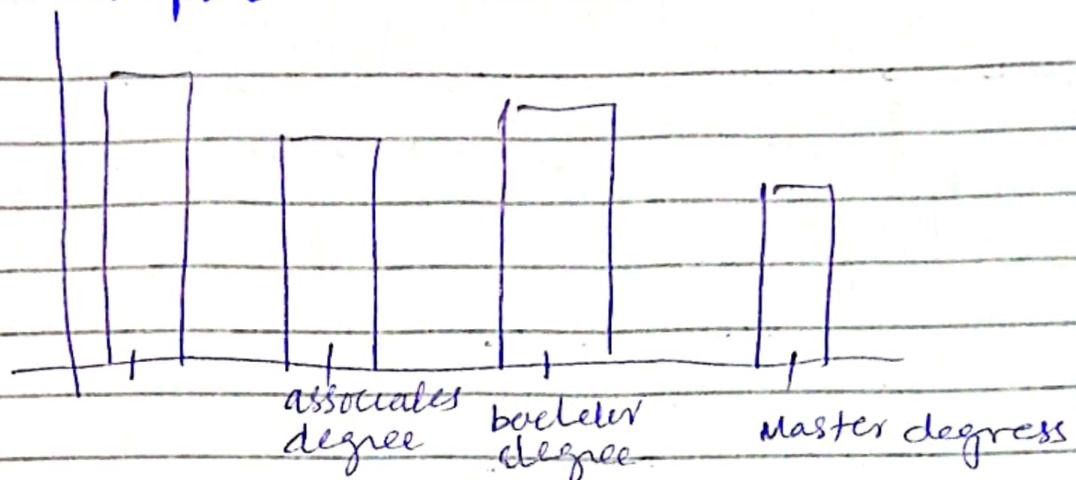
- countplot()

- Counts number of rows per category -

- barplot()

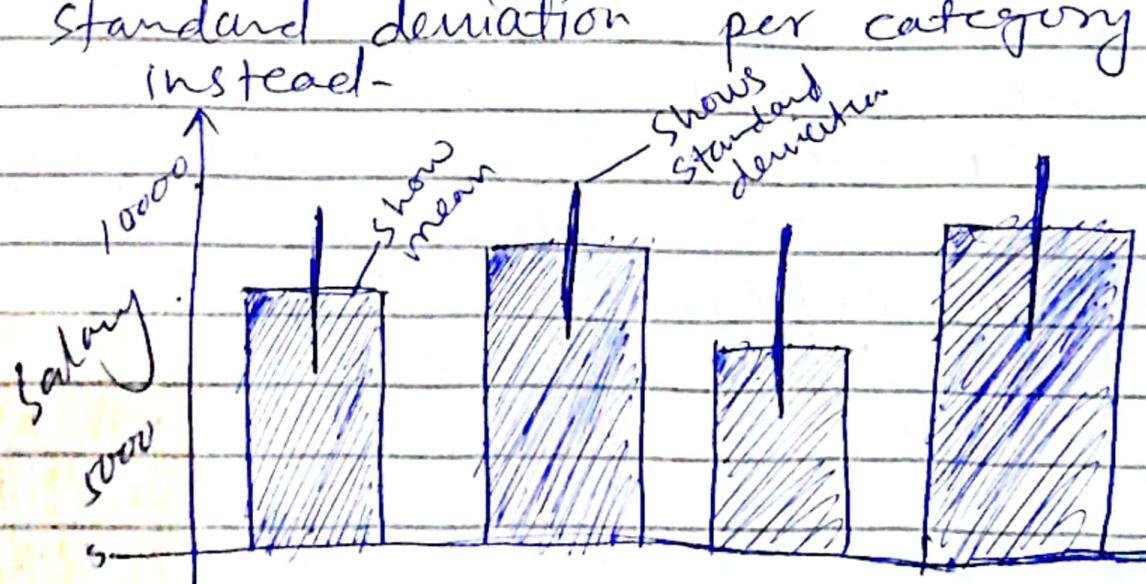
- General form of displaying any chosen metrics per category -

## Count plot:-



## Barplot :-

- The barplot is a general measure which allow us to chose any measure or estimator for the y-axis.
- We could plot the mean value and Standard deviation per category instead-



## 6-Categorical plot :- part two

Statistical Estimation within Categories

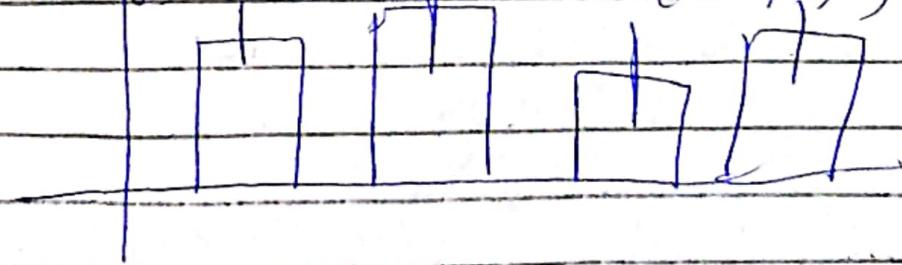
part two: Coding the plots

df.head()

```
df['division'].value_counts()  
sns.countplot(data=df, x='division')  
plt.xlim(100) - not good
```

```
sns.countplot(data=df, x='level of edn', hue='division')  
↓ palette = 'Set2'
```

```
sns.barplot(data=df, x='level of education',  
y='salary', estimator=np.mean,  
ci='sd')  
hue='division'  
plt.legend(bbox_to_anchor=(1.05, 1))
```



## 007- Categorical plots -

### Distribution within Categories

part = One

#### # Distribution within Categories

- Boxplot
- Violinplot
- Swarmplot
- Boxenplot (letter-value point)
- Let explore understanding these plots on the previous salary datasets.

### Boxplot:-

The boxplot displays distribution of a continuous variable -

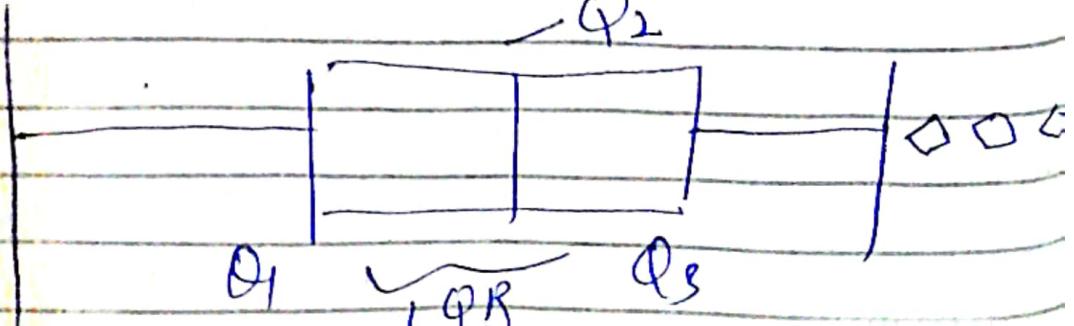
- It does this through the use of quartiles -

◦ Quartiles separate out the data into 4 ~~points~~ equal points - number

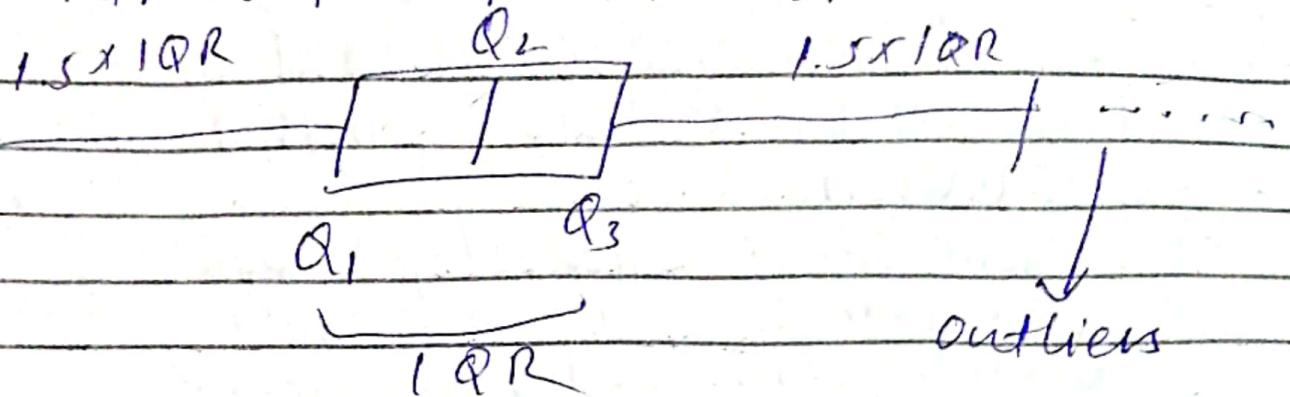
◦ 25% of data points are in bottom quartile -

◦ 50<sup>th</sup> percent ( $Q_2$ ) is the median

$Q_2$



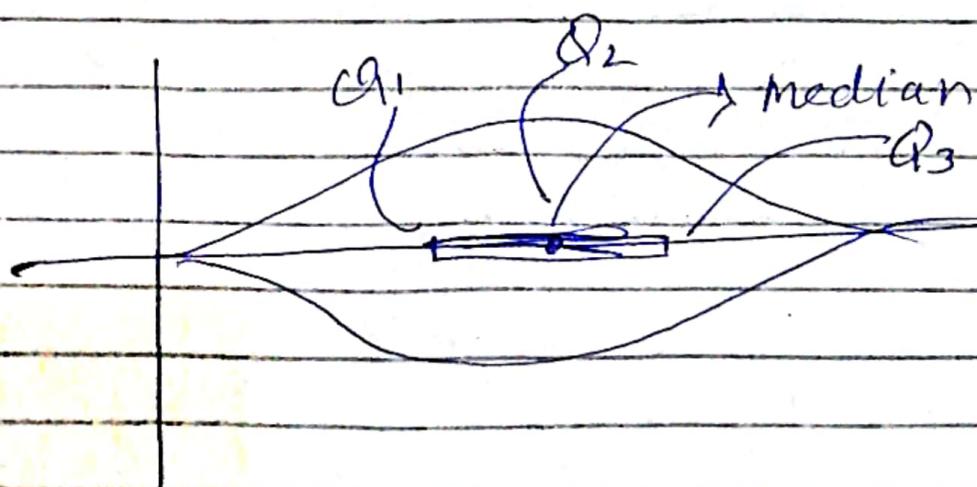
- Median is 50th percentile
- Median splits data in half
- ~~IQR defines the box width~~



## ~~Violin~~ Violin plot :-

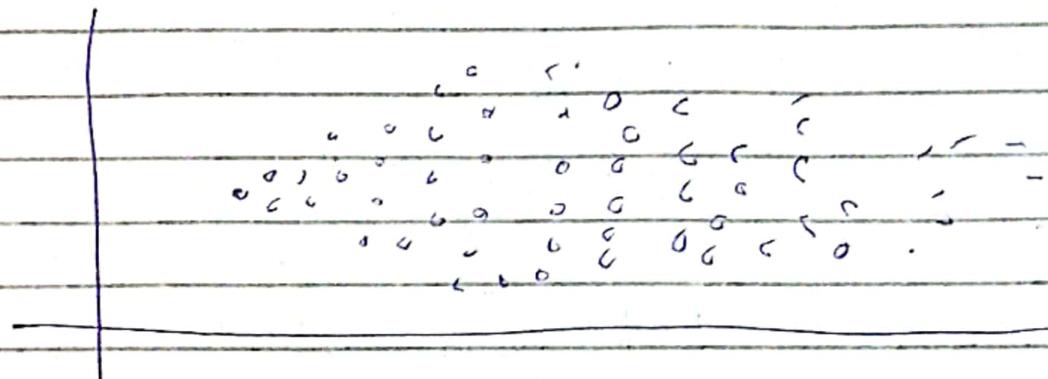
- The violin plot plays a similar role as box plot -
- It displays the probability density across the data using KDE.
- We can imagine it as a mirror KDE

Plot -



# Swarmplot

- Swarmplot is very simple and simply shows all the data points in the distribution.
- For very large data sets, it won't show all the points, but will display the general distribution of them.



will shows how many people we have per category

# boxenplot (letter-value plot)

boxenplot is a relatively new plot developed in 2011 by Claus Heike Hofmann, Karen Leafadaar and Hadley Wickham.

- It's mainly designed as an expansion upon the normal box plot.
- Make sure to read the linked paper in the notebook if you end up using this plot!

## 8 - Categorical plot - Distribution within categories .

### Coding \*

```
df = pd.read_csv('StudentsPerformance.csv')  
df.head()
```

```
sns.boxplot(data=df, y='math score')
```

```
sns.boxplot(data=df, y='math score',  
            x='test preparation course')
```

x = 'parental level'

hue = 'test prepare'

```
plt.legend(bbox_to_anchor=(1, 0.5, 0.4))
```

```
violinplot()
```

```
sns.violinplot(data=df, x='reading Score',  
                y='parental level of education',  
                hue='test preparation course',  
                palette='Set2')
```

split = True,

inner = None

inner = 'quartile'

quartiles = True

Remove inner  
box plot -

$bw =$  — bandwidth stdev in  
Gaussian noise  
(use smaller value)  
(can not good plot larger values)

`plt.figure(figsize=(8,4), dpi=200)`

`sns.swarmplot (data=df, x='mathscore',  
y='gender', size=2,  
hue='test prep course',  
dodge=True)`

dodge  
`dodge = True`  
(split the  
hue to two graph)

`sns.boxenplot(x="math score", y="test prep course",  
data=df, hue='gender')`

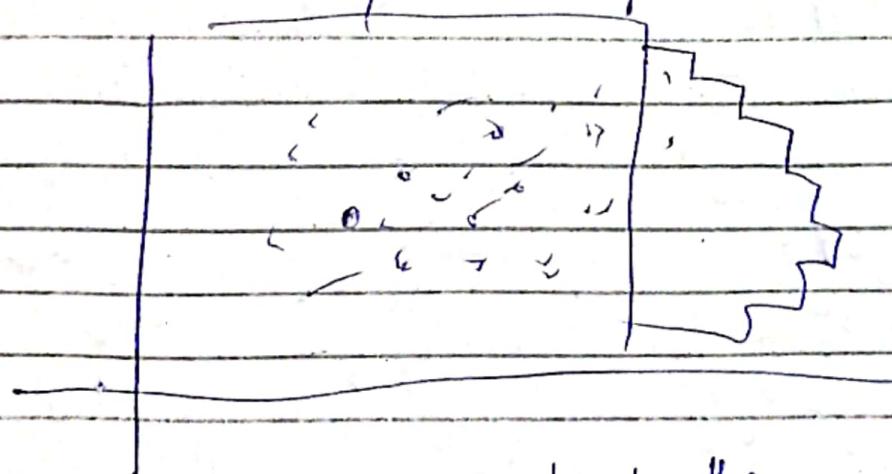
## 9. Comparison plots

part one = Understanding the plots-

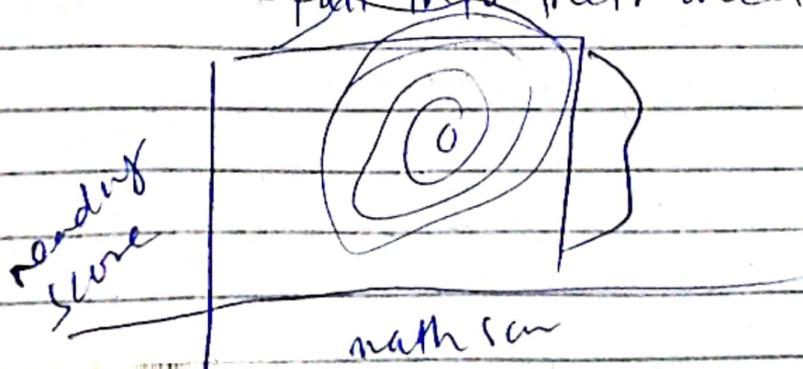
- Comparison plots are essentially 2-dimensional versions of the plots we've learned about so far
- The two mainplots types discussed here:
  - jointplot()
  - pairplot()

jointplot :-

- We can map histogram to each feature of a scatterplot to clarify the distribution within each feature.
- We can also adjust the scatterplot to be hex plot or a 2D KDE plot.



Hexagons are dark the more points fall into their area

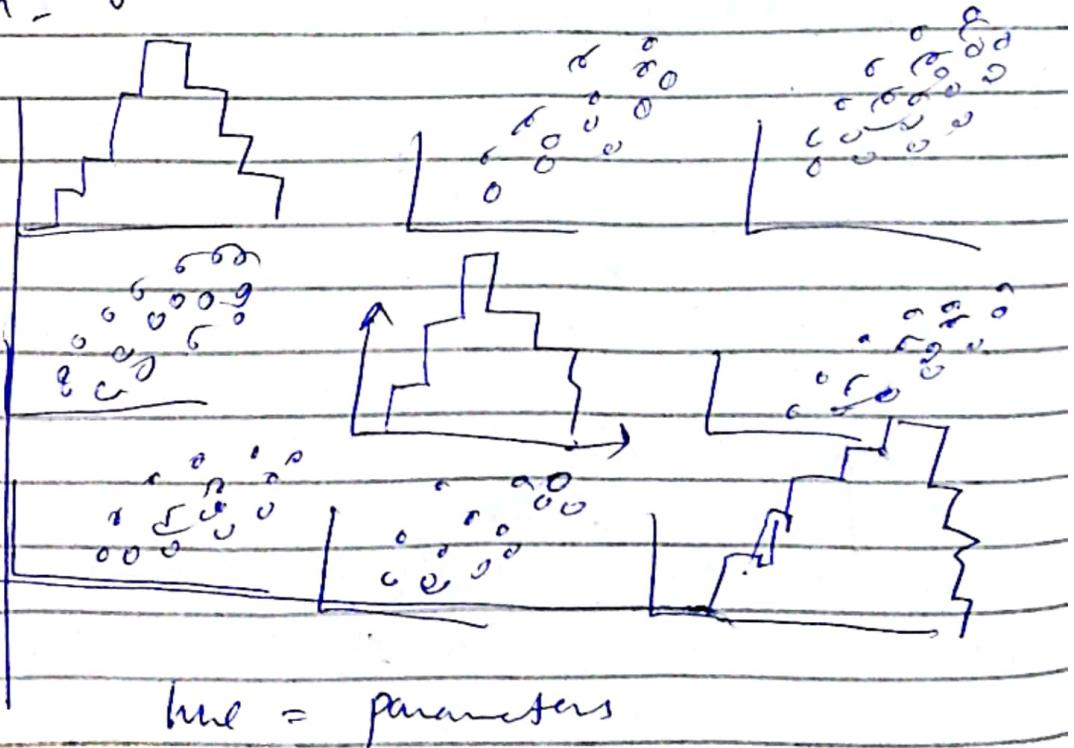


## o pairplot()

- o The pairplot is superuseful. Is a quick way to compare all the numeric columns in a DataFrame
- o It automatically creates a histogram for each column and a scatterplot comparison between all possible combination of columns.

Warning!

- pairplot() can be CPU and RAM intensive for large DataFrame with many columns.
- Its good idea to first filter down to only the columns you are interested in.



## 10. Comparison plots

part two = coding the plots

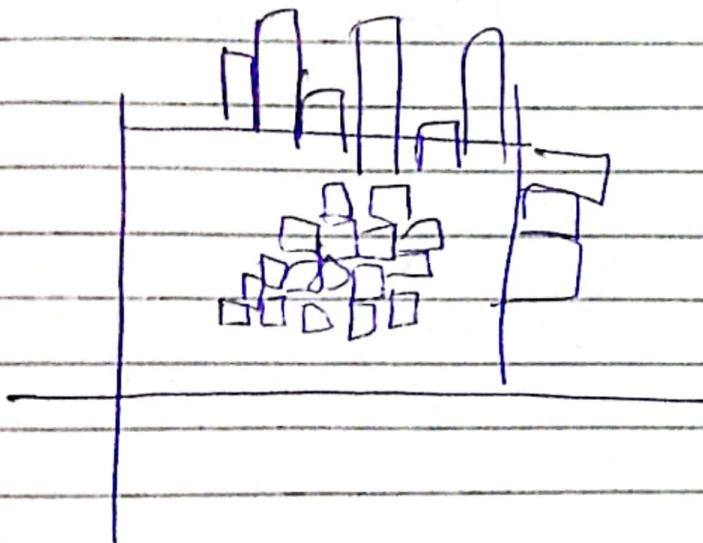
df.head()

sns.jointplot(data=df, x='math score',  
y='reading score')

kind = 'scatter'

kind = 'hex'

alpha = 0.2



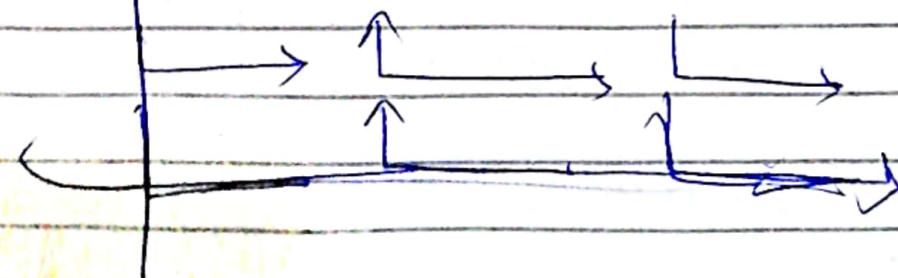
kind = 'hist'

kind = 'kde'

shade = True

hue = 'gender')

sns.pairplot(data=df), hue = 'gender',  
diag\_kind = 'hist')



```
sns.pairplot(data=df, hue='gender', corner  
=True)
```