# EDA
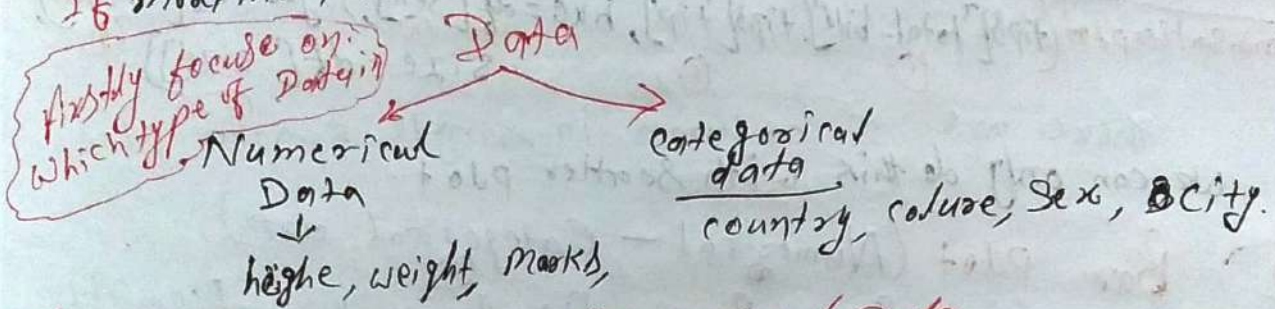
## Univariate Analysis



Single Column → If you Calculate (Analysis) on Single Column it's call Single ~~Variate~~ Univariate Analysis.

If Analysis on two columns → Biy Variate Analysis.

If morel then two columns → Multy Variate Analysis.

Firstly focuse on: Which type of Data in → Data

Numerical Data
↓
height, weight, marks,

Categorical data
country, coluze, Sex, & City.

### Step ① ~~focuse~~ work on Categorical Data.
ⓘ count plot, PieChart!

ⓘⓘ df['Pclass'].value_counts().plot(~~kid~~ kind = 'Pie',
autopct = '%.2f')

### Step ⑪ Numerical Data:
(a) → Histogram → Histogram is a tools for checking the Range of the data

Eg →

| | | |
|---|---|---|
| 10 — 20 Age | → | 15 peoples |
| 20 — 30 Age | → | 40 P.s |
| 30 — 40 Age | → | 30 peoples. |
| 40 — 50 Age | → | 20 peoples. |

↳ Range (Bins)

matplotlib

plt.hist(df['Age'])

You can know about Skewness from this plot.

Same

KDE

df['Age'].skew()
L → Left
R → Right (+)

(b) Distplot (seaborn)
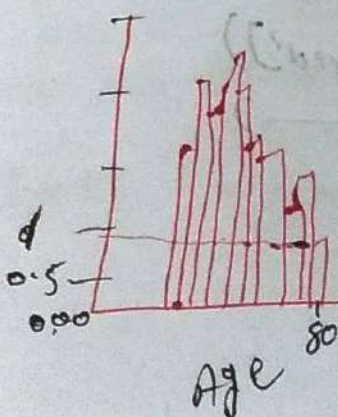sns.distplot(df['Age'])

Kde → kernal density estimation.

Es plot se hum probability k bare me pta chlta hai

Q → Age 80 hone ka Probability kya hai

Ans → 1% chanse hai Age 80 hone ka.

© Box plot.
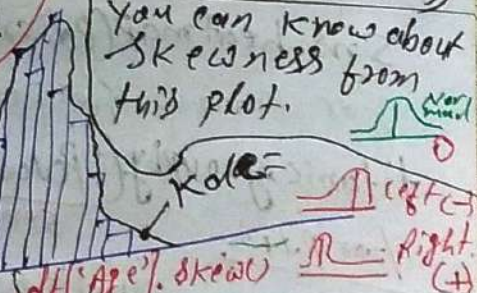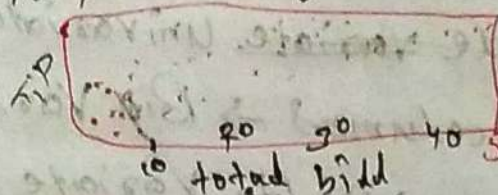sns.boxplot(df['Age'])



Age 80

Age

# EDA Bivariate

Work on 2 columns

① Scatter plot (Numerical-Numerical Data)
     check Relation

(tips dataset)

.sns.scatterplot (tips['total-bill']. tips['tips']) hue =df['sex]



Said-jaise bill badh rha
hai waise-waise tip V badh
rha hai,

## Multivariate

sns.Satterplot(tips['total-bill']. tips['tip'], hue =df['sex'], style =df['smoker],
                        ①               ①                ③                        , size =df['size']) ④
                                                                                            ⑤

there are 5 column in single scatter plot
we can only do this with scatter plot.

② Bar plot (Numerical — Categorical data)
                          ①                         ⑪

Sns.barplot(titanic['Pclass'], titanic['Age'], hue =titanic['sex'])

③ Box plot (Numerical- categorical)

Sns.boxplot(titanic['Sex'], titanic['Age'], hue = titanic['Survived'])

④ Distplot (Numerical- categorical)

sns.distplot(titanic[titanic['Survived']==0]['Age'], hist =false)
sns._____ " _____ "1 _____ " _____ )

⑤ Heat map (Categorical-categorical)

sns.heatmap(pd.crosstab(titanic['Pclass'], titanic['Survived']))

for know the percentage

titanic.groupby('Pclass'). mean()['Survived']×100 → p% show
                                      sen

for plot

(titanic.groupby('Pclass'). mean['Survived']× 100). plot(kind ='bar')
                                       sex

⑥ ClusterMap (Categorical - Categorical)

pd.crosstab(titanic['Sibsp'], titanic['Survived'))

for plot sns.clustermap(_____))  _____

⑦ Pair plot
　　sns. pairpdot (iris, hue = 'spacies')
　　↳ it will self decide and plot with same
　　column , (Cate → Cate )
　　　　　　　　name → nume )

⑧ Lineplot ( Numerical - Nameof cud )
　　↳ it is basicly use for date and time dataset
　　　　　　　　　　　　　　　　　　　[Share market graph]

new = flights. groupby ('year'). sum(). reset_index()
　　sns. lineplot (new['year'], new['passengers])

for pivot table

flight. pivot_table (values= 'Passengers', index = 'month', columns='yrar)

Sns. heat map( _____ " _____ )
sns. clustermap( _____ " _____ )