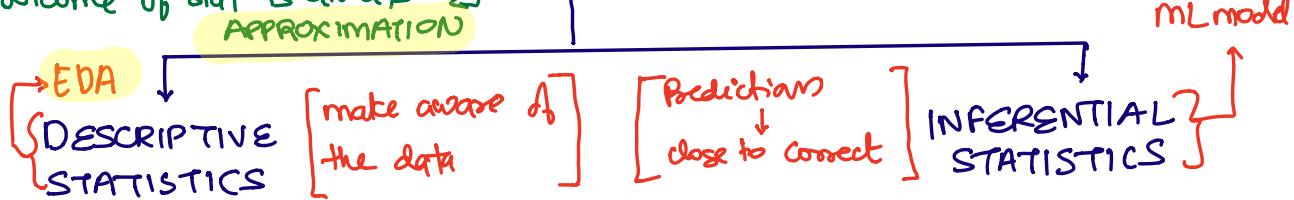


working with a sample ←
 outcome of stat is always ←
 APPROXIMATION

Statistics

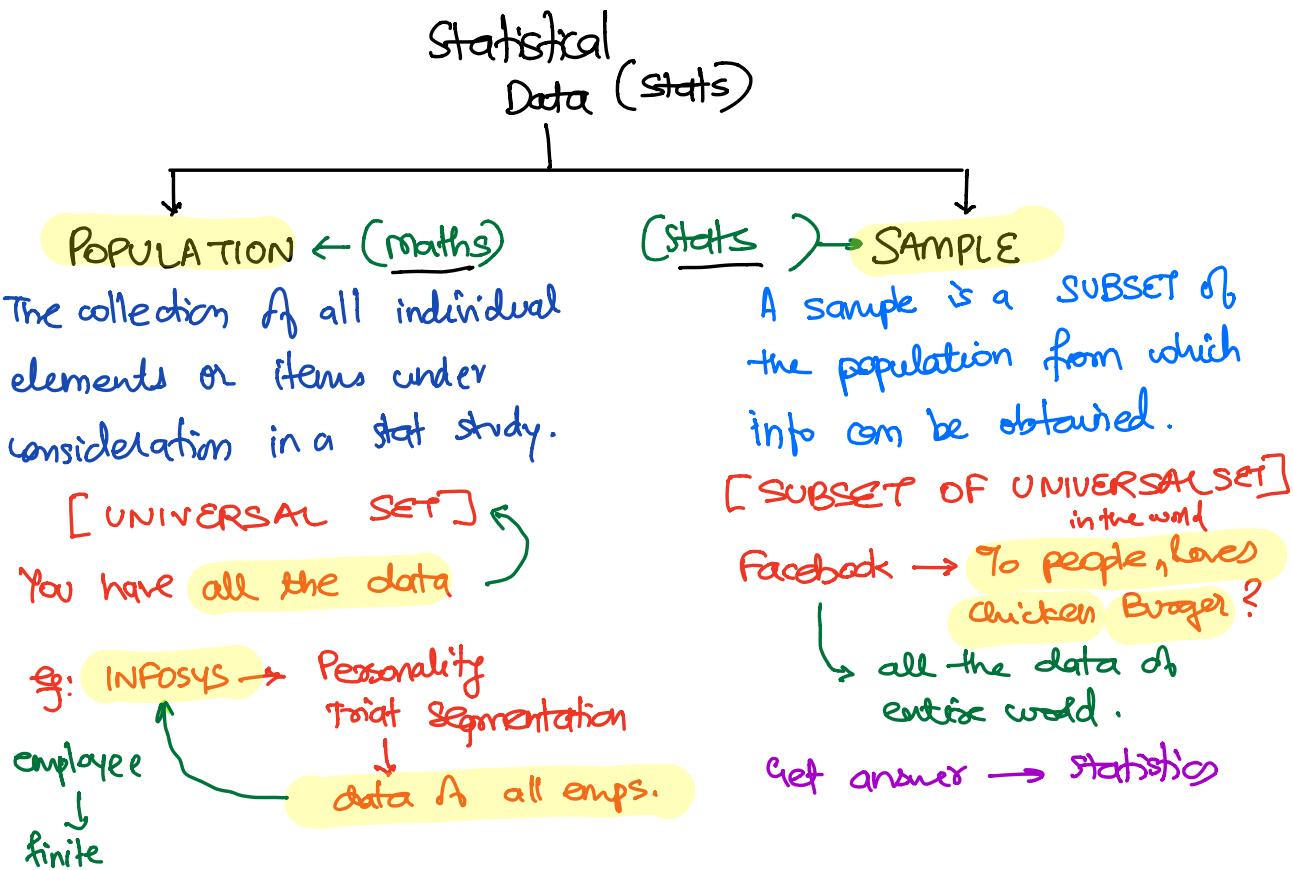
→ a way to get information
 out the data.



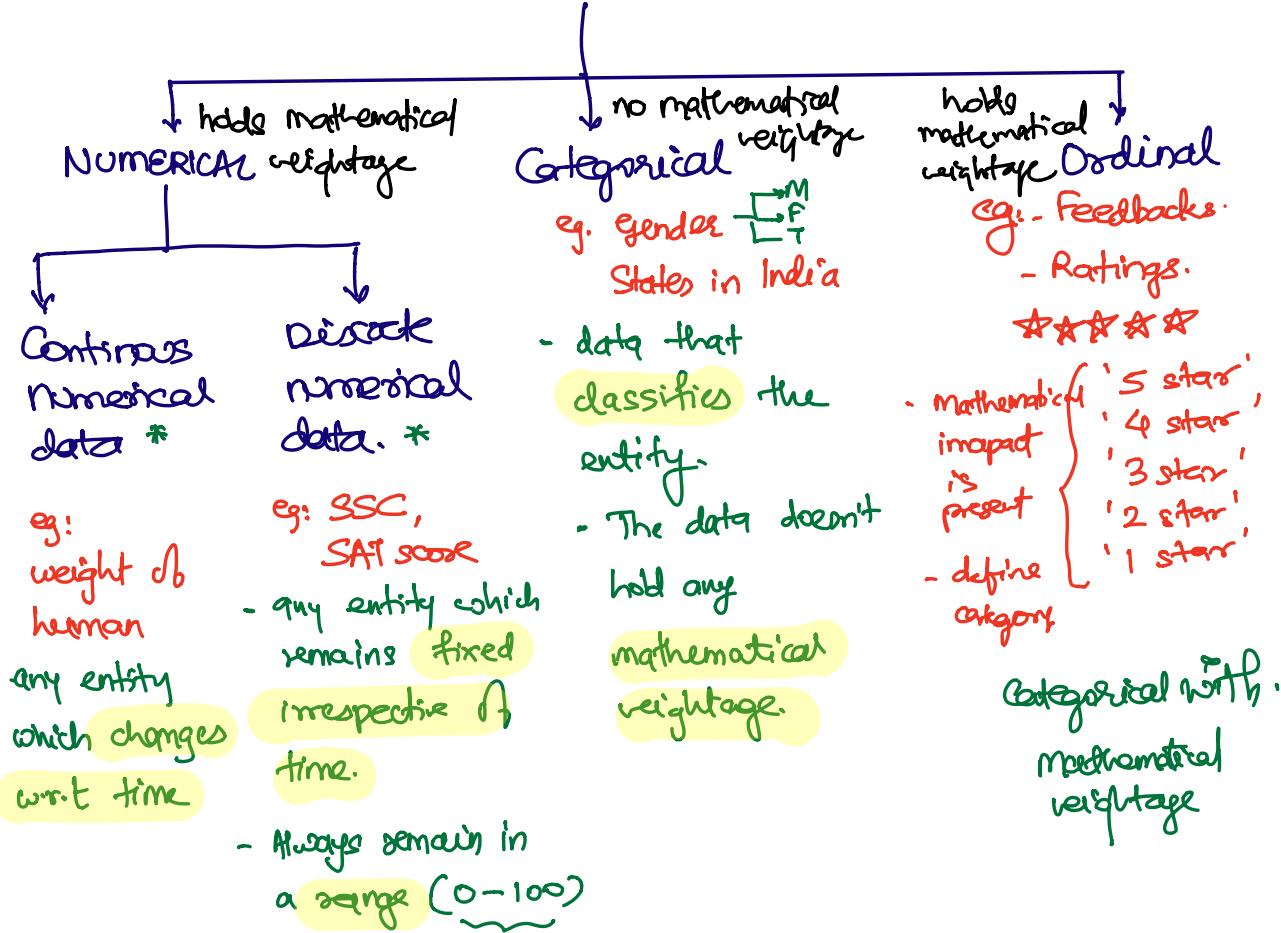
Technique for organising
 and summarizing info.
 using

1. Table
2. Graph
3. Charts

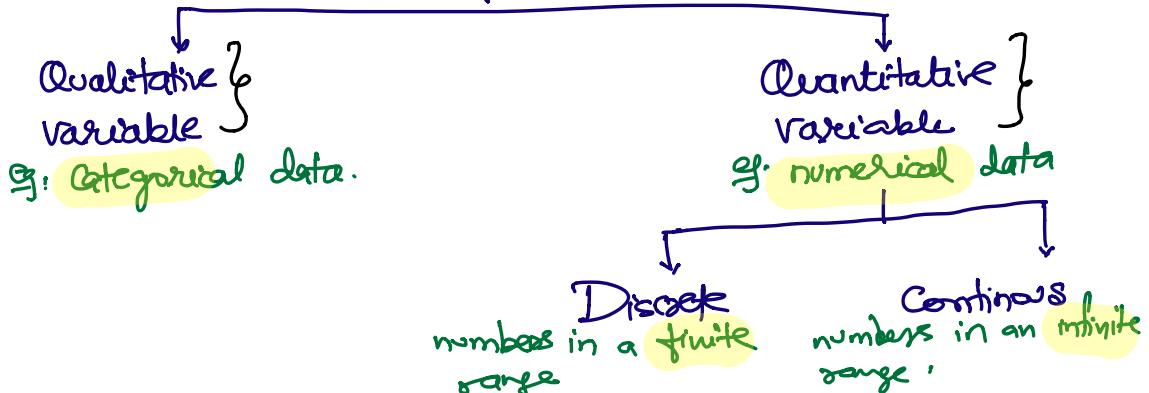
Technique for drawing and measuring the reliability of conclusions of the population
 based on information.



Types of data you deal with in Data Science



Types of Variables

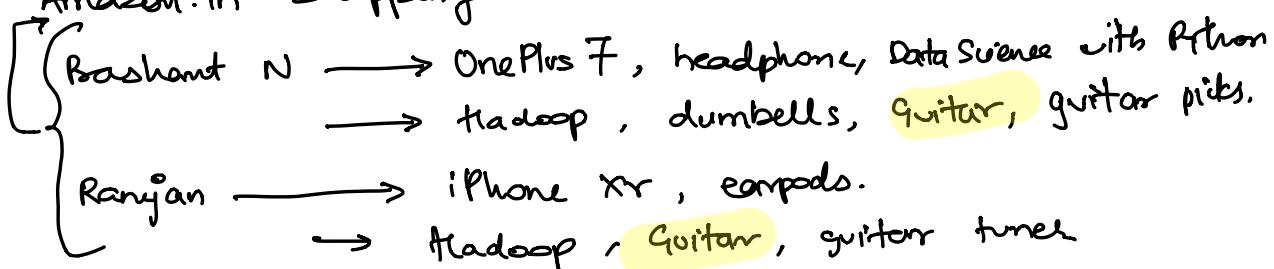


DESCRIPTIVE STATISTICS

[Data Exploration Phase]
Exploratory Data Analysis
(EDA)

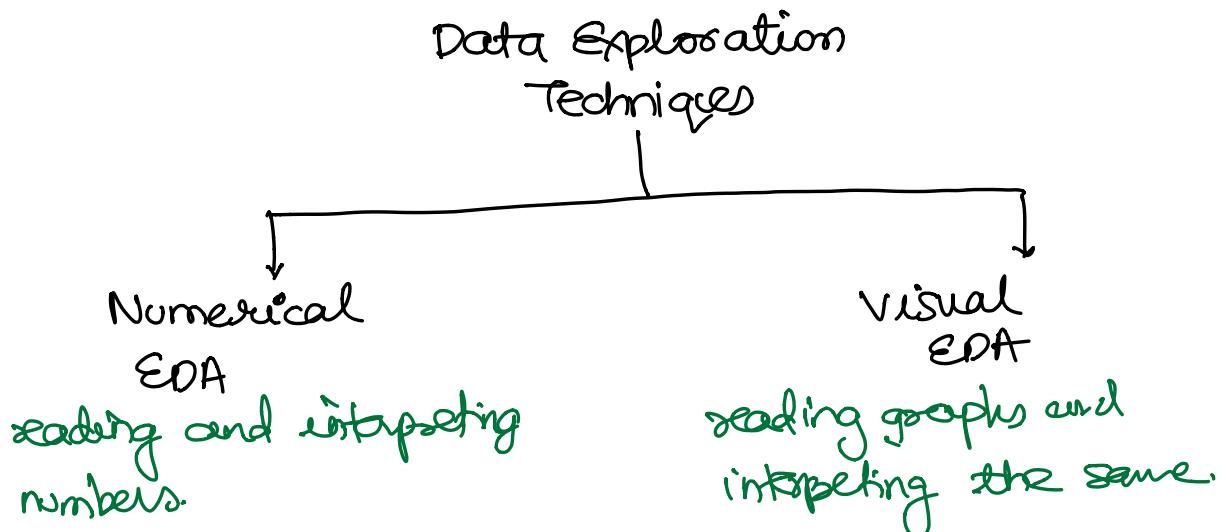
- Descriptive stats is the first step to deal with the sample
- mindset of a Data Scientist while performing **Descriptive Stats** include:
 - To check the distribution of data
 - To identify and remove OUTLIERS
 - To identify and deal with INAPPROPRIATE data.
 - To identify and deal with MISSING values in data.
 - To help detect any kind of ASSOCIATION / ^(Association Rule mining)
 \nearrow RELATIONSHIP / \nearrow PATTERN between two variable
in a dataset \uparrow

Amazon.in Shopping website.



Nikhil → Hadoop
Guitar

Amazon will recommend
Nikhil



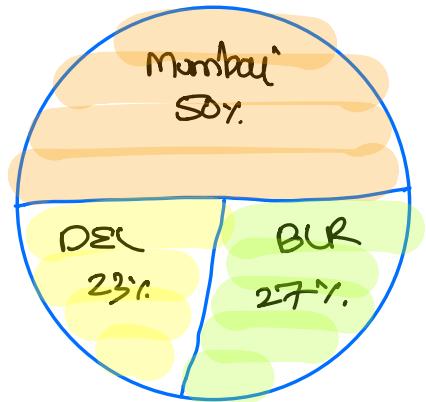
Visual EDA

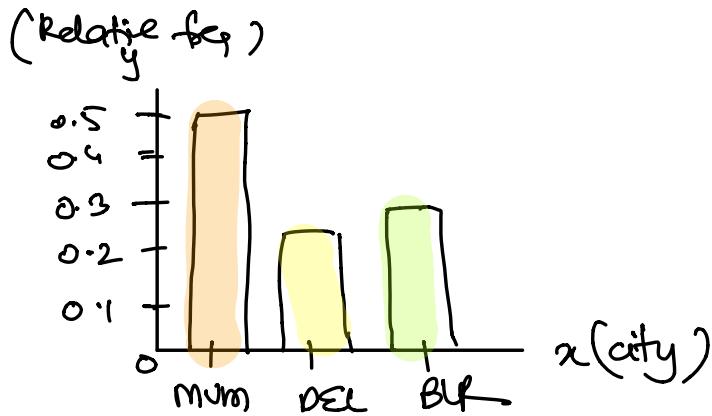
- ① Dealing with Qualitative data. (**Categorical data**)
Represent qualitative data using two approach

Frequency Distribution Table

city	Tally	freq.	R.F
MUM		9	9/18 → 0.5 / 50%
DEL		4	4/18 → 0.23 / 23%
BLR		5	5/18 → 0.27 / 27%
		18	

Graph → **Pie**
Bar



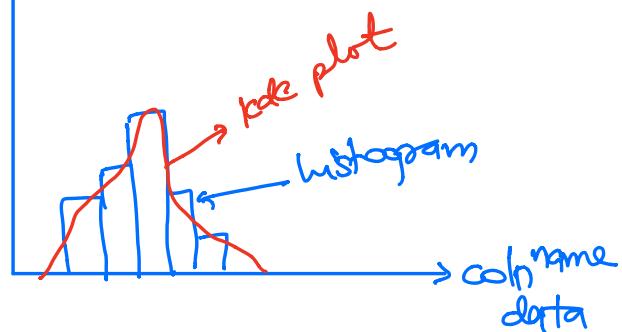


Quantitative Data

freq dist table

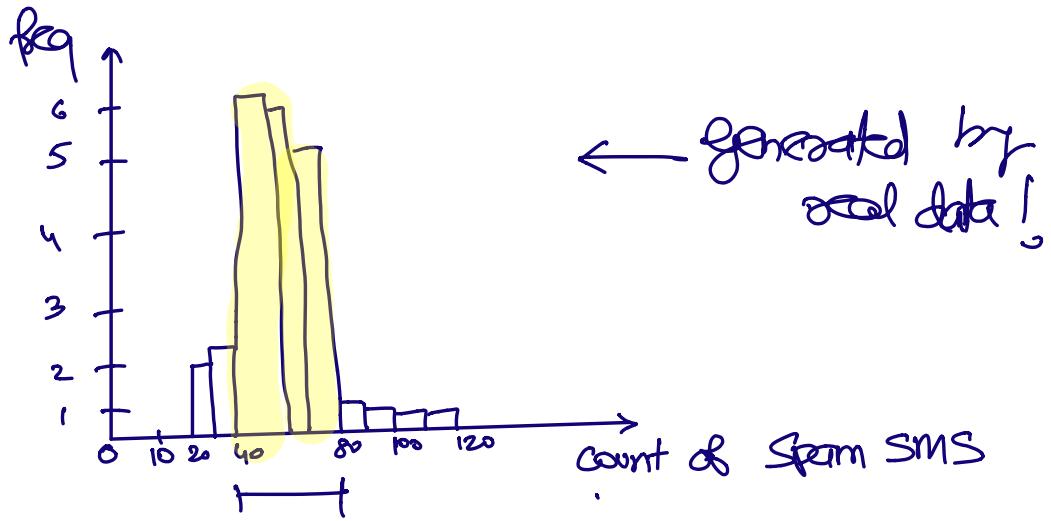
(numerical data)

histograms / KDE plot
freq. ↗

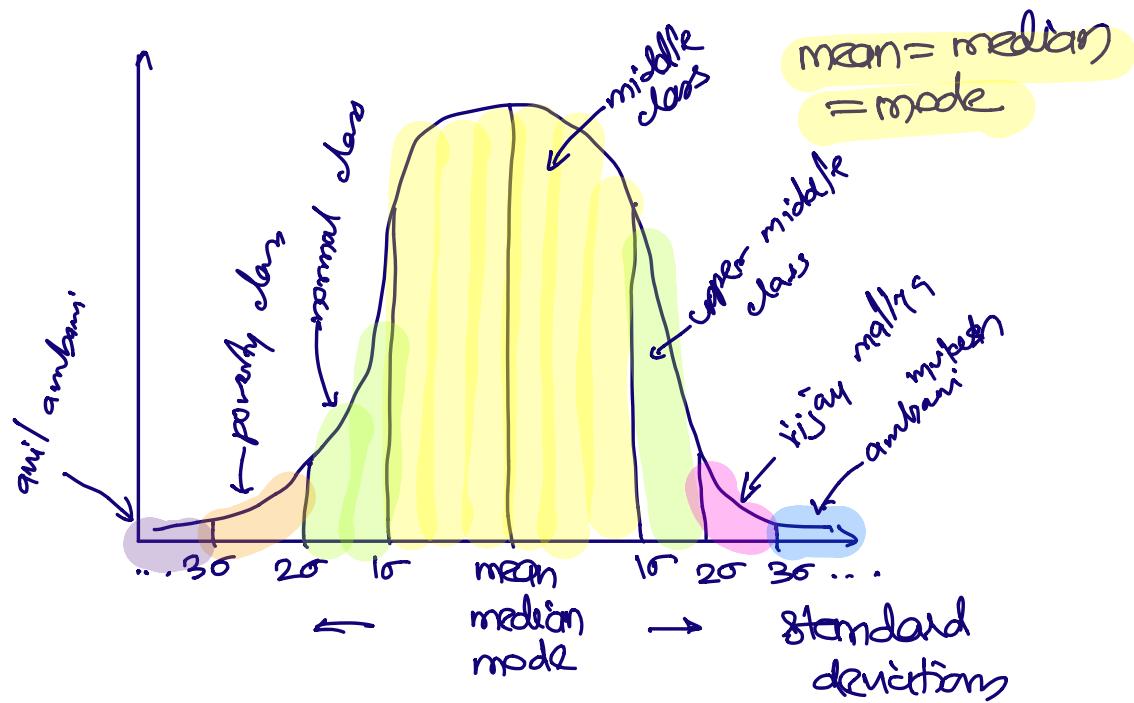


DATA DISTRIBUTION → only applicable for
(defining the partial nature of the data)

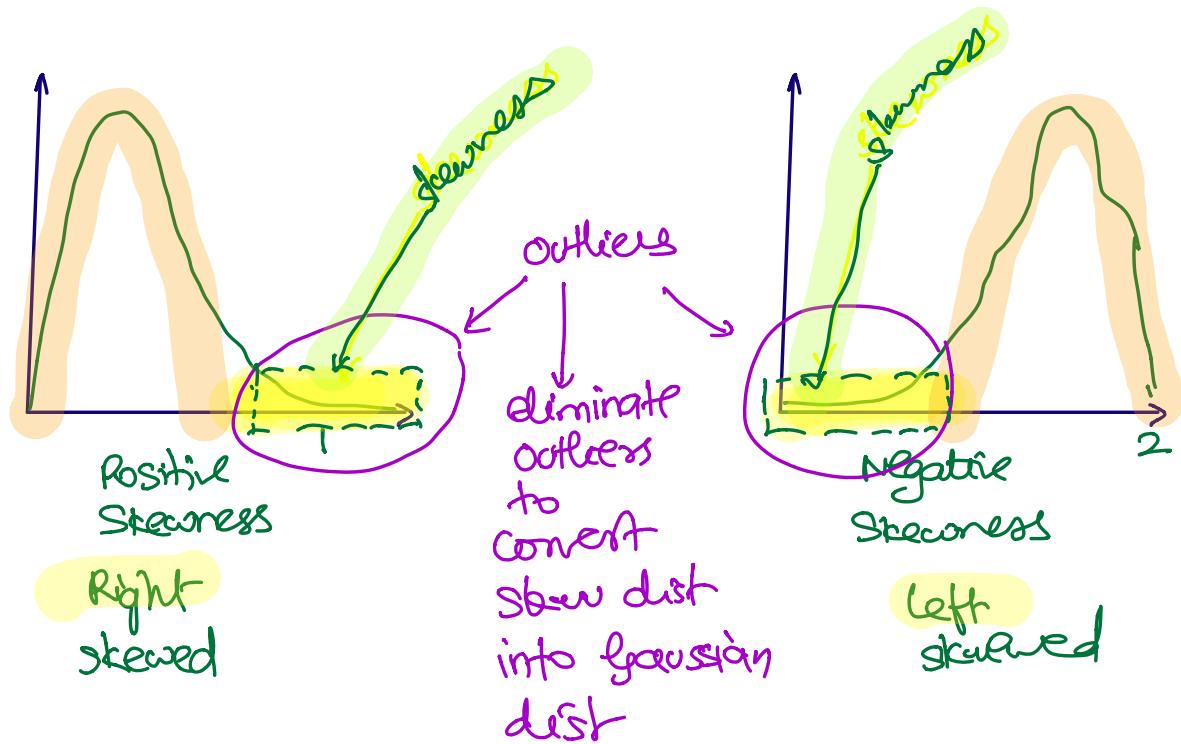
Quantitative data!
(Numerical → continuous
Discrete)



② Gaussian distribution (Normal distribution)

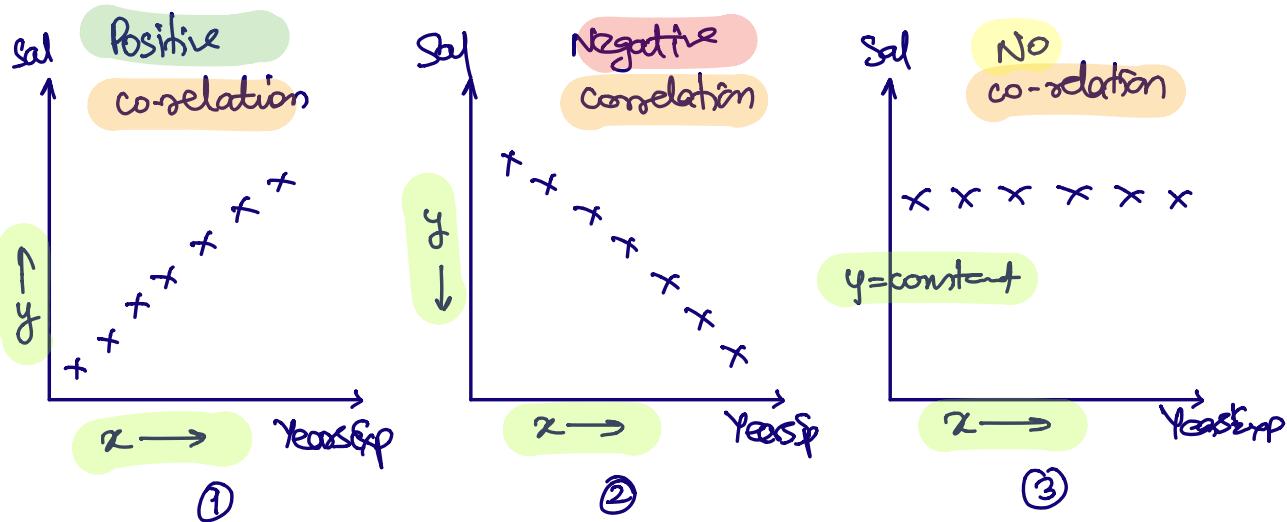


Skew distribution



Relationships in data distribution

linear relationship. ← Present } Regression analysis



OUTLIERS → Outliers are extreme values that affect

Stat perspective @ Gaussian nature of data.

Domain perspective B) Domain nature of the data.

e.g. Indian Stock market → NSE (High volatility)

market falling / Crashing ↓

AXIS Bank → Buying

To detect outliers and same the same using STAT perspective, you can use the concept of Quartiles.

Quartiles

min

Q1 25%

Q2 50% ← mean | median | mode (Posely subject to dist. of data)

Q3 75%

Q4 max

IQR (remove outlier columnwise)

Inter Quartile Range → $Q_3 - Q_1$

$$\text{lower Range} = Q_1 - (1.5 * \text{IQR})$$

$$\text{upper Range} = Q_3 + (1.5 * \text{IQR})$$

Algo :-

① Ensure your column data is sorted in ascending order

② Get value of Q_3 and Q_1 and calc IQR.
using the below,

$$\text{IQR} = Q_3 - Q_1$$

③ Calc. lower range = $Q_1 - (1.5 * \text{IQR})$

④ Calc. upper range = $Q_3 + (1.5 * \text{IQR})$

