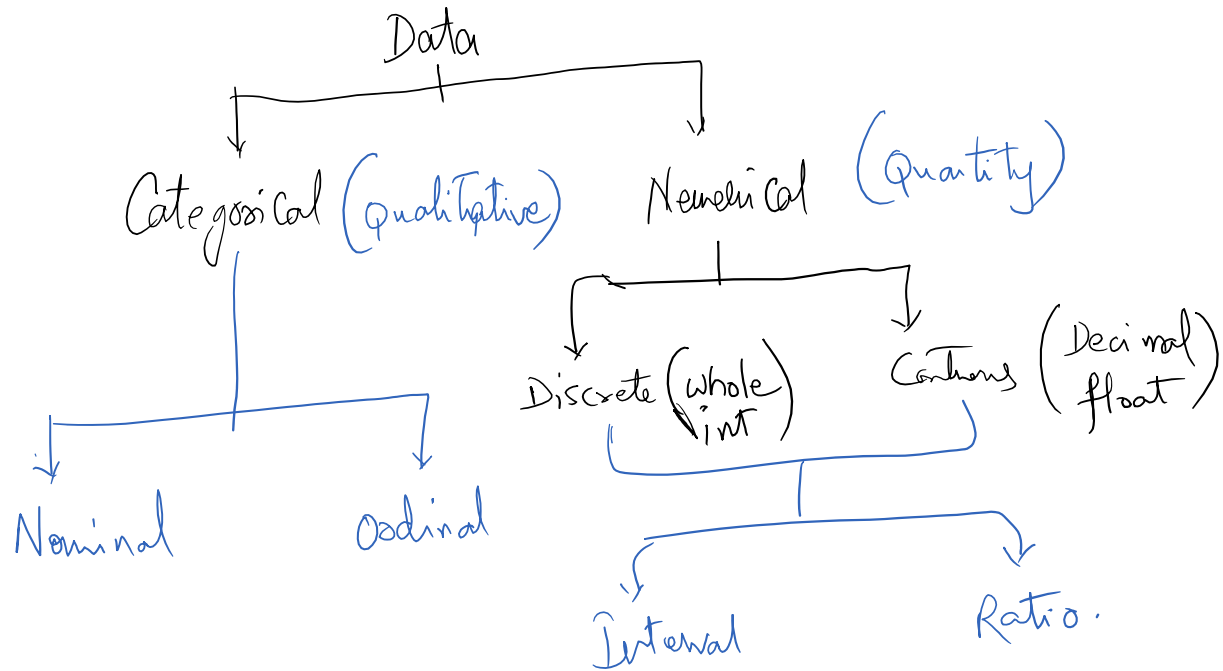
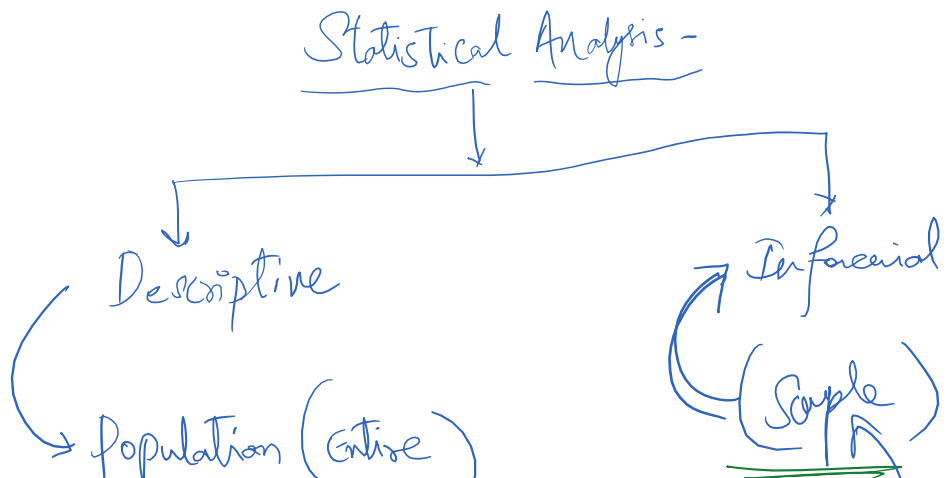
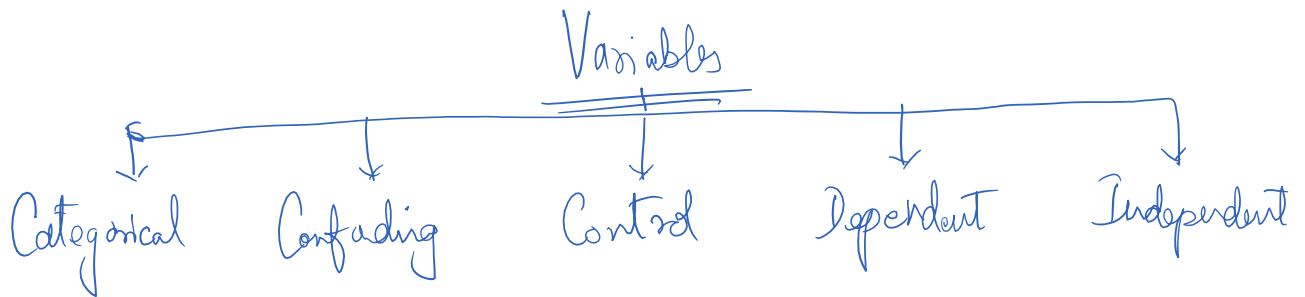


Data → set of Recorded facts → Information : processed data with meaning

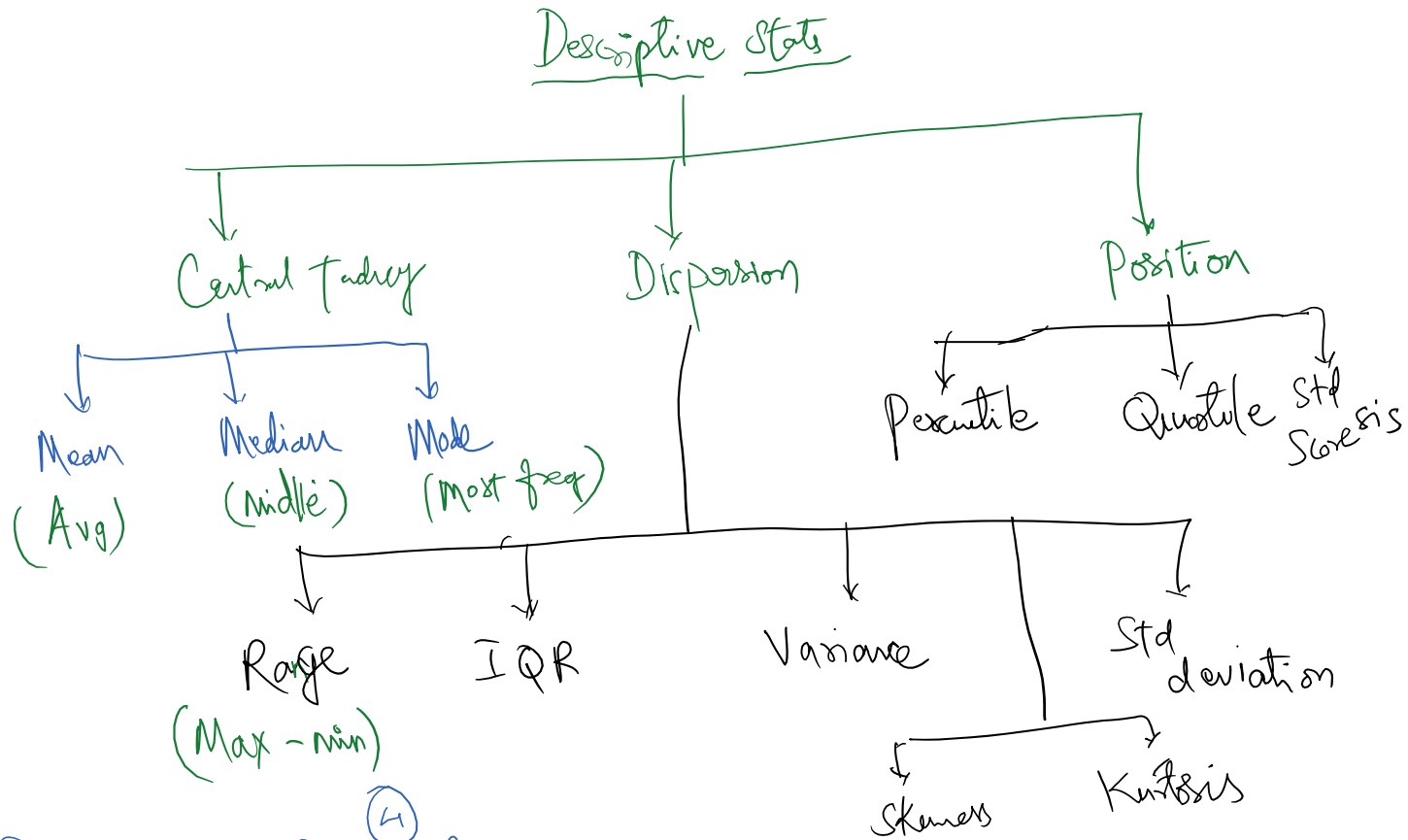


Variables : Represents an unknown value / value that changes

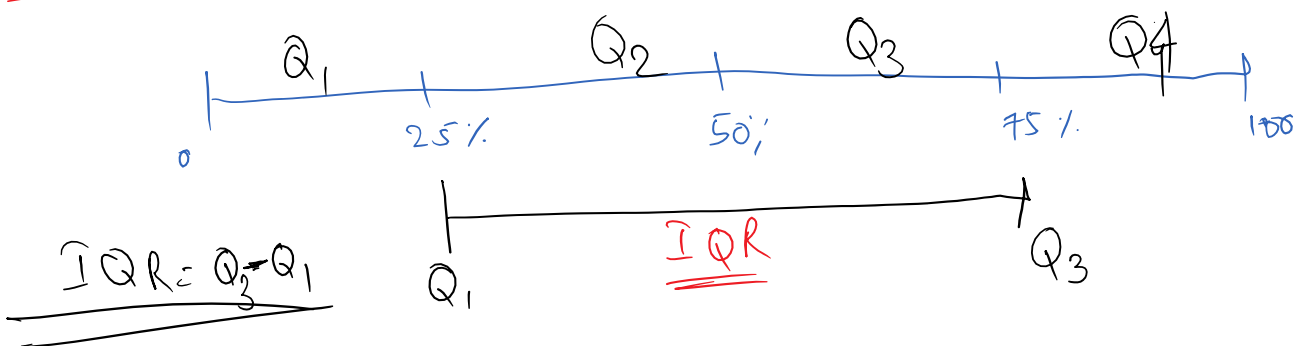


Population (entire data)

Sample  
From Population



IQR → Inter Quartile Range <sup>(4)</sup>



Variance → Measurement of how far a number is from mean.

$$S^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

$x$  → data point

$\bar{x}$  → mean

$n$  → Total datapts.

$$S^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

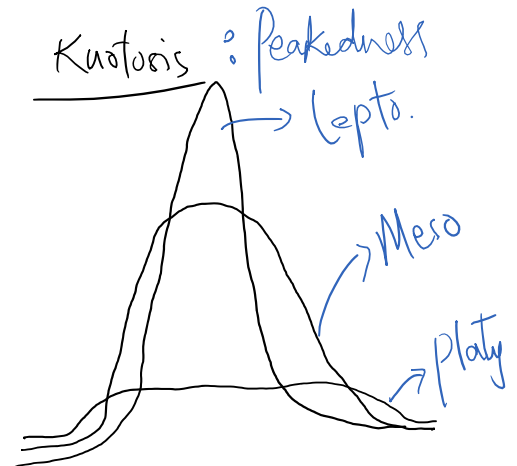
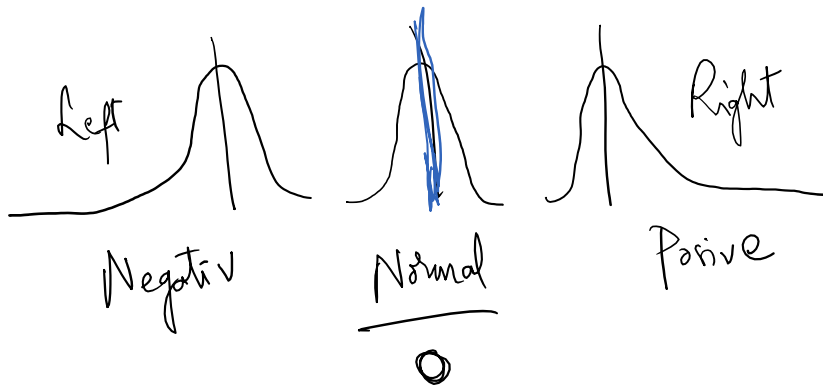
$\bar{x} \rightarrow$  mean

$n \rightarrow$  Total datapts.

Std dev  $\rightarrow$  Root of Variance

$$\sigma = \sqrt{\text{Variance}}$$

Skewness : Symmetry



Percentile : Divide data in to 100 parts equally  $\rightarrow$  each part is 1 percentile

Quartile :  $\frac{1}{4}$  parts.  $\rightarrow$  Quartile

Standard Score : How many standard deviation, a value is away from Mean

Z Score

$$Z = \frac{x - \bar{x}}{\sigma}$$

$x \rightarrow$  Data point

$\bar{x} \rightarrow$  Mean

$\sigma \rightarrow$  Std deviation

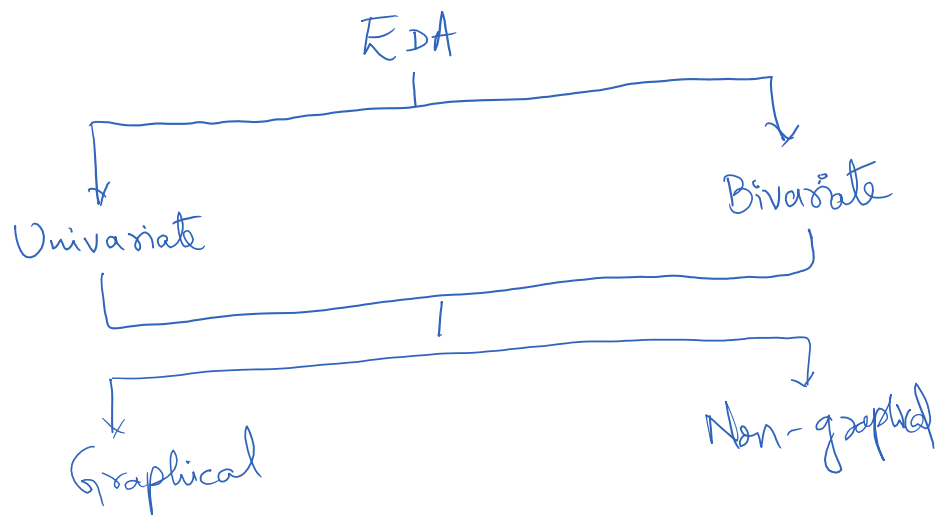
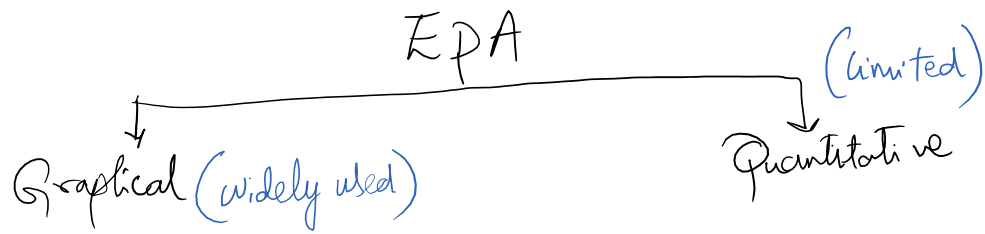
Age  $\rightarrow$  Salary  $\rightarrow$  Different Scale

23	25K
24	26K
32	50K
46	75K

Age(z)	Sal(z)
0.3	0.8
0.4	0.85
0.6	0.93
0.8	0.96

EDA → Exploratory data Analysis.

- Insight to the data
- Extract important Variables
- Detect outliers & Anomalies.
- Understand data better



Univariate → Non-graphical EDA

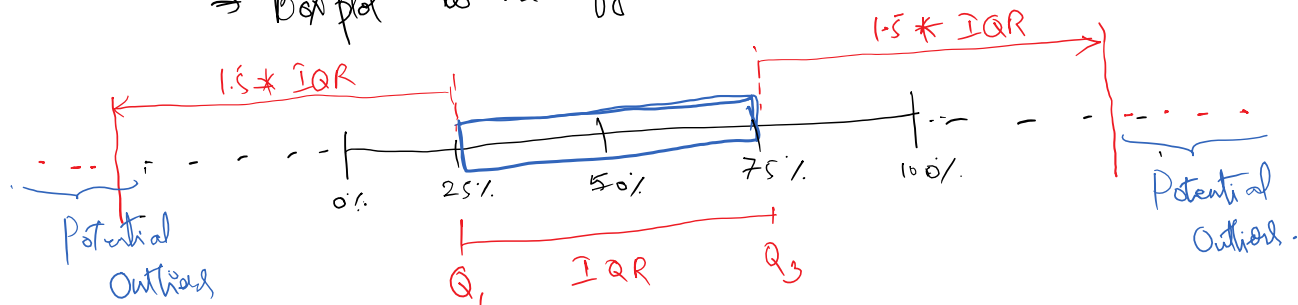
- ⇒ Measure of Centre (Mean, Median, Mode)
- ⇒ Measure of Dispersion (Range, IQR, Variance, Stddev)
- ⇒ Measure of Position (percentile, Quartile, Std dev)
- ⇒ Measure of Shape (Skewness (-, 0, +), Kurtosis)
- ⇒ Test for Outliers

Outlier : Data that deviates highly from main stream data.

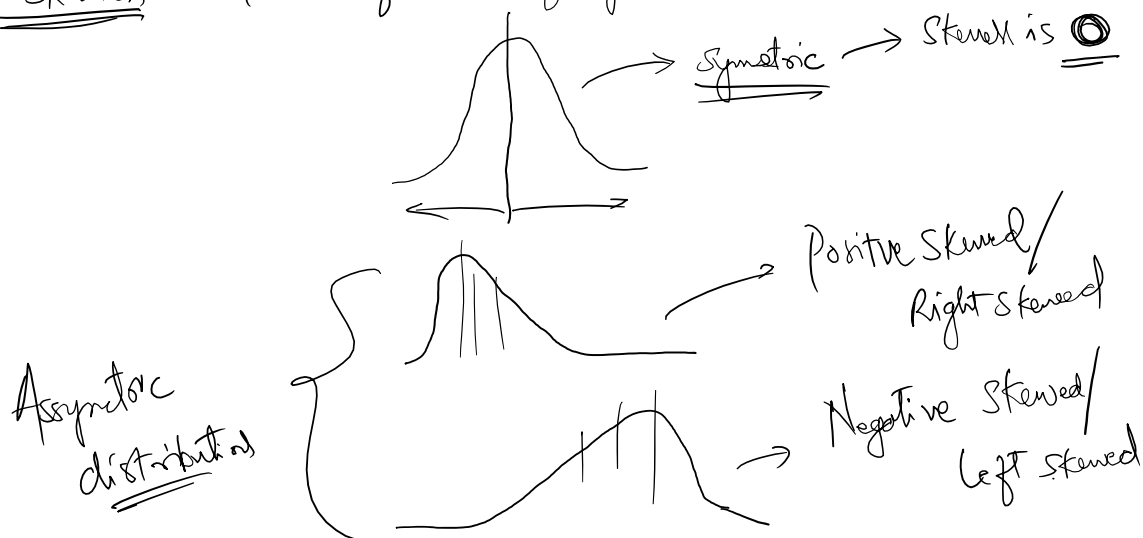
Outlier : Data that deviates highly from main stream data.

Data  $> / < 1.5 * \underline{IQR}$   $\rightarrow$  Generally considered outlier

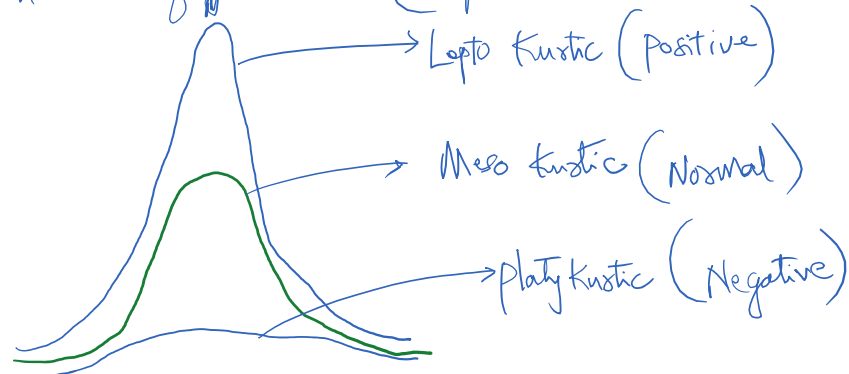
$\rightarrow$  Box plot to identify outliers.



Skewness : Measure of lack of Symmetry



Kurtosis : Measure of Peakedness (Compared to Normal dist)



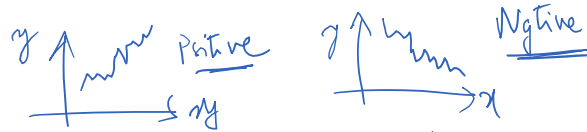
Univariate Graphic EDA  $\rightarrow$  Shape & Distribution

$\rightarrow$  Bar graph

- Histograms
- Pie chart
- Box plot → Quartiles → Outliers
- Scatter plot

Bivariate EDA: Plot the graphical relation b/w Variables.

Bivariate EDA → Non graphic.



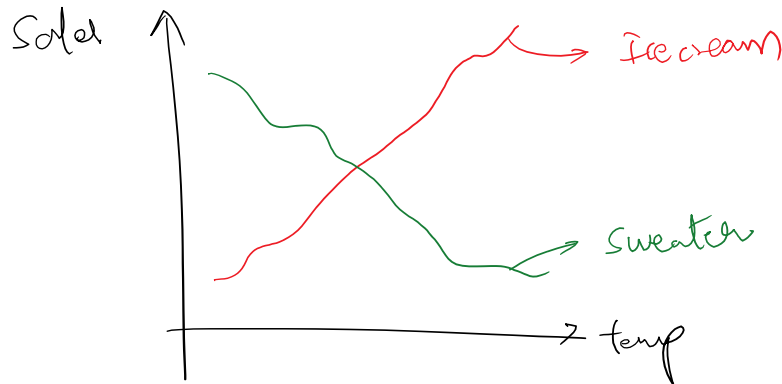
- Covariance
  - Correlation
- Help to identify the linear relation b/w Variables.

Covariance

- Var's are +/− related
- Only direction

Correlation (Association)

- Degree of the relation of Variables
- Direction + Straight



Correlation → Pearson, Spearman, Kendall

	Var1	Var2	Var3
Var1	1	0.3	0.8
Var2	0.2	1	0.4

Var 2	0.3	1	-0.4
Var 3	0.8	-0.4	1

Covariation : Change in one variable is similar to another which has no relation with the variable

Bivariate EDA : Graphical

- Scatter plot (Outliers, Positive/Neg, Density, Trend)
- Heat Map (Highlight the important & not important one)

Guide for graphic EDA

Variable

Categorical

Univariate + Continuous

Bivariate + Continuous

GEA

Bar, pie

Line plot, Histogram

Scatterplot

Objective

Idea of Distribution

Outliers

Relation b/w Variables

Graphic EDA

Histogram

Box, Scatter, Histogram

Scatter, Cov & Corr

Cross tab

State	Gender	People	Total
-------	--------	--------	-------

State	Gender	People	Total
KA	M	2	7
	F	5	
TN	M	3	9
	F	6	
KL	M	1	3
	F	2	
AP	M	4	12
	F	8	