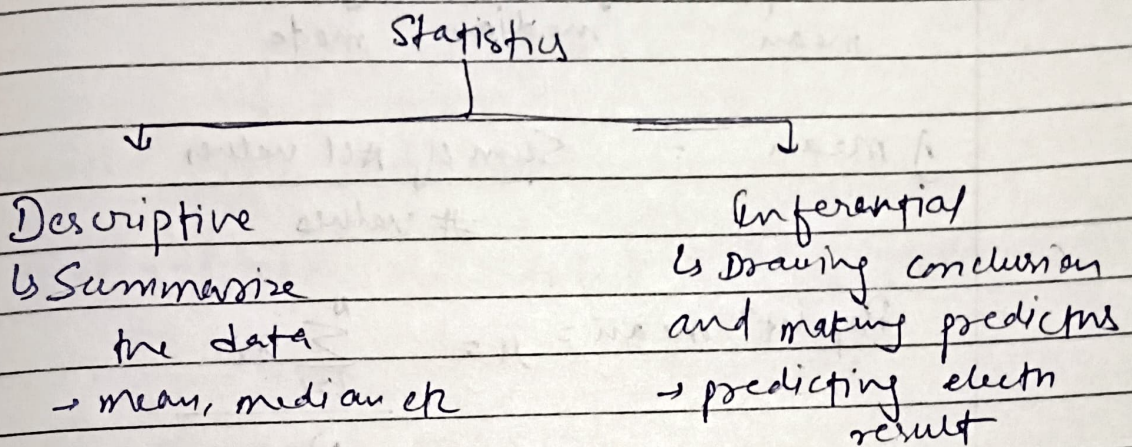


Session-38

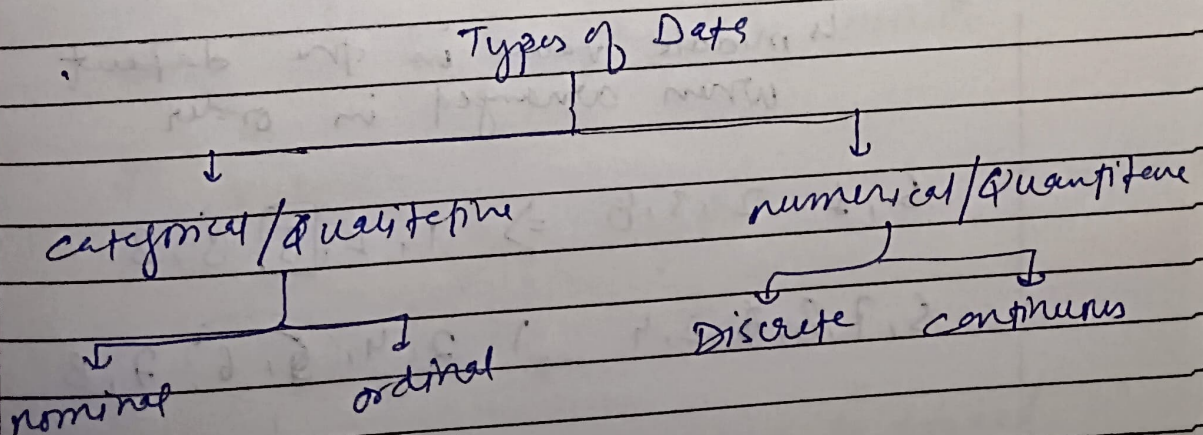
Descriptive Statistics - I

Statistics Diagram ✓



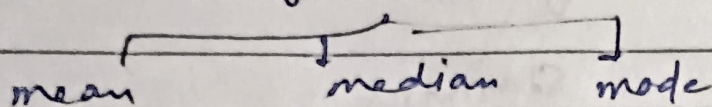
Population v/s Sample ✓

- Inferential Stats -
 - Hypothesis testing
 - Confidence interval
 - ANOVA
 - Regression analysis.
 - Chi-Square test
 - Sampling Techniques



• Central tendency -

↳ Single value that represent center of data



1) mean = $\frac{\text{Sum of All values}}{\# \text{ values}}$

Population mean = $\mu = \frac{\sum_{i=1}^N x_i}{N}$

Sample mean = $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$

$N \rightarrow$ No. of items in population
 $n \rightarrow$ " " " " Sample

1) ~~Median~~ : If outlier, then we may get present wrong mean.

2) median -

↳ middle value in the dataset when arranged in order

5, 1, 2, 3, 8 \Rightarrow 1, 2, 3, 5, 8

5, 7, 8, 8, 2, 4 \Rightarrow 2, 4, 5, 8, 7, 8

↓
 $\frac{5+8}{2} = 5.5$

3) mode - most frequent value in Dataset

1, 1, 5, 1, 2, 3, 2, 5. \therefore mode = 1

4) weighted mean -

	1	2	3
weight	0.2	0.3	0.8

$$\text{weighted mean} = \frac{1 \times 0.2 + 2 \times 0.3 + 3 \times 0.8}{0.2 + 0.3 + 0.8}$$

5) Trimmed mean

1, 2, 3 } 8, 10, 12 } 20, 22
 trim it out calculate its mean trim it out

• measure of Dispersion -

↳ how much spread data is

1) Range \rightarrow (min. value, max. value)

↳ It is affected by outlier

2) Variance -

$$\sigma^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}$$

• Variance is proportional to spread.

↙ Can't be used in inferential stats

1. Mean Absolute Deviation

$$= \frac{\sum |x_i - \bar{x}|}{n}$$

↳ less prone to outlier

• Variance is highly prone to outlier.

• Population variance = $\frac{\sum (x - \mu)^2}{N}$
 σ^2

• Sample variance = $\frac{\sum (x - \bar{x})^2}{n-1}$
 s^2

3) Standard deviation -
 $\sqrt{\text{variance}}$

• Its unit is same as Data.

4) Coefficient of Variation (CV)

$$CV = \left(\frac{\sigma}{\mu} \right) * 100$$

• Salary and age can't be compared.

By calculating CV - we can tell which of Salary & age are more closer to their mean.

• Graphs for univariate Analysis

↳ Depends whether data is categorical or numerical

1) Categorical column -

i) frequency distributⁿ table

category	frequency

⇒ using this we can build Bar graph

ii) Relative frequency

↳ we ~~get~~ get % of frequency

↳ we can build pie chart

iii) Cumulative frequency -

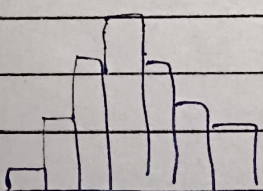
A	fy	c. freq
1	20	20
4	30	50
2	40	90

2) Numerical column

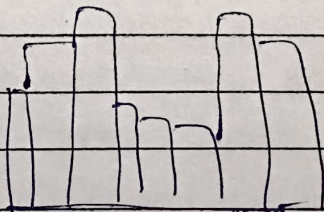
↳ frequency Distributⁿ ✓

↳ Histogram ✓

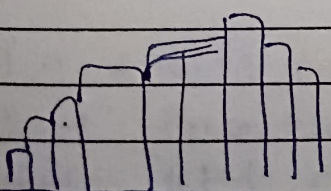
↳ we use bins



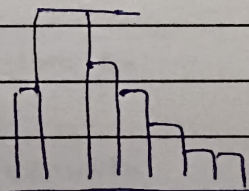
Symmetric



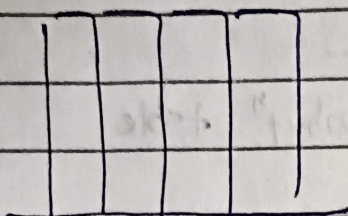
Bimodal



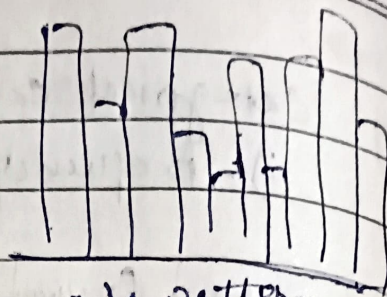
left skew



Right skew



uniform



No pattern.

• Graphs for Bivariate analysis -

↳ cat - cat

↳ cat - numerical

↳ Num - Num

1) Categorical - categorical

↳ we make contingency table A / cross table

2) Numerical - Numerical

↳ Scatter plot

3) categorical - Numerical

↳ a lot of thing can be done

↳ can make contingency.