

Statistics

↳ The science from learning the data.

Descriptive stats → summarizing data

Inferrential stats → methods for drawing conclusions about a population based on info from a sample.

Population

Entire Collection

Sample

Subcollection of members selected from a population.

Step I → Identify the question

Step II → Collect the data

Step III → Analyze the data

Step IV → Draw conclusions

Parameter vs statistics

↓
describes a population

describes a sample

Type of data you dealing

PAGE NO.

DATE

Quantitative

Vs

Qualitative

(Numerical data)
(Counting)

(Categorical data)
(Labeling)

Discrete Data

Continuous Data

(certain values)

(Any value
within range)

Missing Value

→ Delete
or → input values (Replace)

Sampling Methods

i) Simple random sample

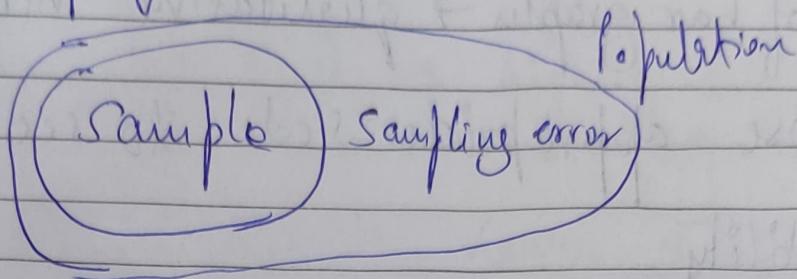
ii) Systematic sample

iii) Stratified sample

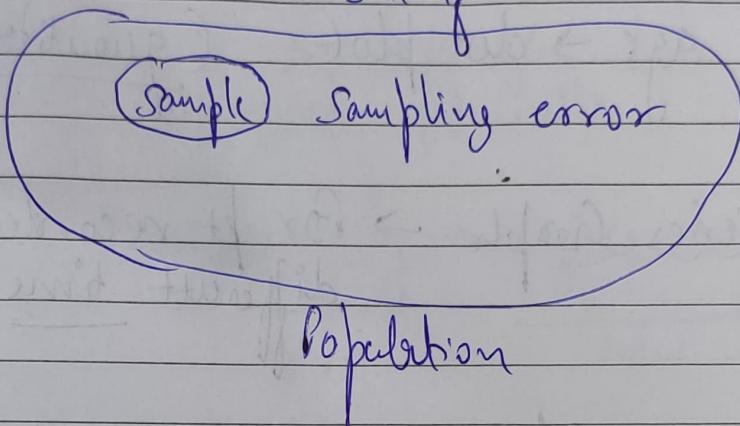
iv) ~~cluster~~ sample
Cluster

v) Convenience Sample

Sampling error



but if



Good Decision of Experience

- i) Randomization
- ii) Replication (Repeat) → wide range
- iii) Blinding (Placebo effects)

Visualization

Use of bar graphs \rightarrow qualitative variable

A Choose appropriate scale to see
volatility

Eg:- age \rightarrow dot plot (quantitative var)

Time Series Graph \rightarrow For ft recorded at different time

Dot plots are good for Comparisons

Histograms \rightarrow can compare well but on same scale.

frequency polygon

\hookrightarrow Histograms dot in middle joined

Classes

\hookrightarrow Dividing data into classes (equal width)

Lower class limit \rightarrow smallest value in each class

Upper class limit \rightarrow highest value in each class

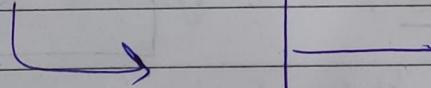
Class mid point $\rightarrow \frac{\text{lower} + \text{upper}}{2}$

Class width \rightarrow diff b/w two consecutive lower limits.

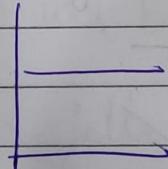
Relative frequency \rightarrow $\frac{\text{frequency for specific class}}{\text{total no of data pts}}$

The column sums up $\frac{100}{\text{ }}$

i) Uniform Distribution \rightarrow Same frequency



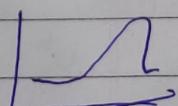
ii) Normal distribution
 $\hookrightarrow \mathcal{N}$



iii) Skewed Distributions

Right distribution \rightarrow Right foot $\rightarrow \mathcal{N}$

Left distribution \rightarrow Left foot $\rightarrow \mathcal{N}$



Paired Data → Related Data

(height & weight)

⇒ Correlation

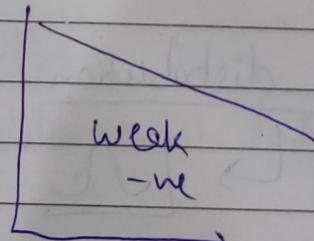
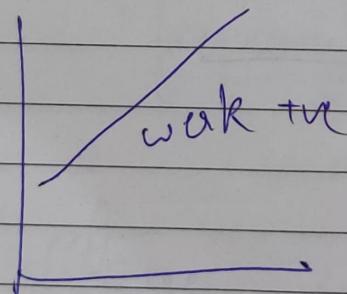
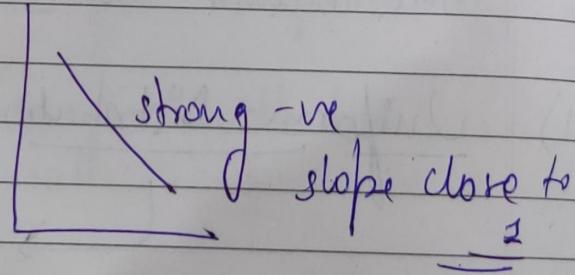
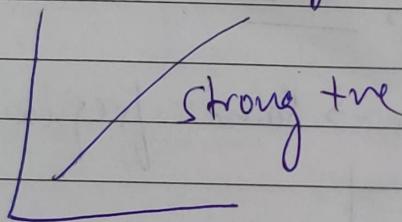
↳ does not imply causation

↳ lurking variable

↳ indirectly involved

Linear correlation → straight line relationship

regression line



Center of your data

Resistance → a measure of center is resistant if the presence of extreme values doesn't cause it to change very much.

Mean → Avg of data → $\frac{\text{sum of all values}}{\text{no of total values}}$

$$\text{Sample} \rightarrow \bar{x} = \frac{\sum x}{n} \rightarrow \text{sample}$$

$$\text{Population Perimeter} \rightarrow \mu = \frac{\sum x}{N} \rightarrow \text{population}$$

So, mean is not resistant

↳ coz any outlier is there, it will be added and hence it will be changed

Median → Value that falls in exact middle of data set → when in a ascending order

So median is resistant to extreme values

Mode :- Greatest frequency of data

The mode is resistant to extreme values
bcz outliers would never be in mode

Symmetric data \rightarrow Mean \sim Median

Measure of Spread

Range \rightarrow Diff b/w largest & smallest
 ↳ Not resistant

Standard Deviation :- measure of how much individual values deviate away from mean.

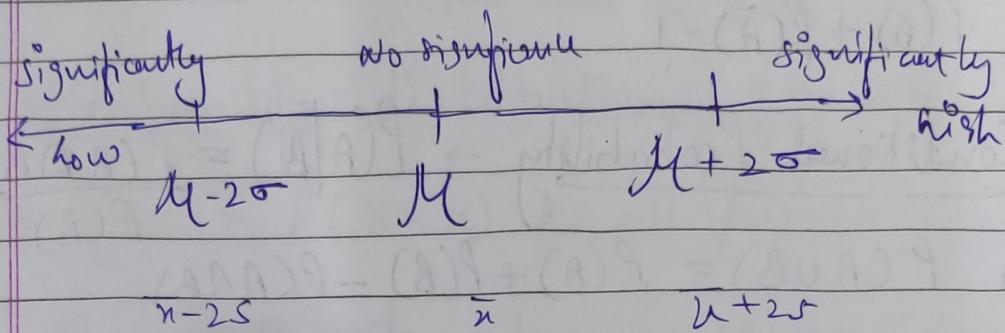
$$\text{Sample} \Rightarrow S = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}} \quad \begin{matrix} \rightarrow \text{like range} \\ \text{from } 1 \text{ to } n \end{matrix}$$

$$\text{Population} \Rightarrow \sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}} \quad \begin{matrix} \text{Degree of freedom} \end{matrix}$$

Prop of SD

- i) Never be negative
- ii) 0 only if all data same
- iii) Large values \rightarrow more variance & wider spread
- iv) Not resistant \rightarrow will have extreme values

$$\text{Variance} = (\text{SD})^2$$



Measure of location

Percentile → divides in 100 groups $\rightarrow 1\% \text{ in each}$

Percentile of x = $\frac{\text{number of values less than } x}{\text{total no of values}} \times 100$

Quartiles $\rightarrow 25\% \text{ of gap}$

$Q_1 \rightarrow 25\%$

$Q_2 \rightarrow 50\% \rightarrow \text{Median}$

$Q_3 \rightarrow 75\%$

where 0% is minimum & 100% is max

$$IQR = Q_3 - Q_1$$

Outlier \rightarrow Above $Q_3 + (1.5 \times IQR)$

Below $Q_1 - (1.5 \times IQR)$

Probability

$$P(A) + P(\bar{A}) = 1$$

Conditional Probability $\rightarrow P(A|B) = \frac{P(A \cap B)}{P(B)}$

$$\Rightarrow P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Mutually exclusive event \rightarrow cannot happen at same time

$$P(A \cap B) = 0$$

\Rightarrow Sampling with replacement vs without
 ↓
 independent not independent

$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(B \cap C) - P(C \cap A) + P(A \cap B \cap C)$$

$$P(A \cap B \cap C) = P(A) \times P(B|A) \times P(C|E \cap B)$$

Contingency table \rightarrow (Venn diagrams)
 ↳ Group histograms table

Discrete random Variable

↳ that take countable values

Probability distribution of n → stable $x=0, 1, \dots$

pmf $\rightarrow f(x) = P(X=x)$ for all possible value of n

Mean discrete value $\rightarrow \mu$ or $E(x)$

$$\sum x \cdot f(x)$$

→ pmf of x

Variance of Discrete random variable

$$\hookrightarrow \text{Var}[X] = \sigma^2 = E[X^2] - (E[X])^2$$

$$\sigma^2 = E[X^2] - \mu^2$$

$$\sum x^2 \cdot f(x)$$

S.D of Discrete random Var

$$SD[X] = \sigma = \sqrt{E[X^2] - (E[X])^2}$$

Bernoulli Experiment \rightarrow Either Success or Fail

|| no \rightarrow sum of all success

Binomial distribution

$$f(x) = P(X=x) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$\text{Mean} = np$$

$$x=0, 1, \dots, n$$

$$\text{Var} = np(1-p)$$