

ALL ABOUT STATISTICS

Lets denote Population as (N) and sample (n)

What all are the sampling techniques, describe about it?

→ a) Simple random sampling :- Every member of a Population (N) has an equal chance of being selected for your sample (n)
eg:- Exit Poll, survey

b) Stratified Sampling :-

Strata \rightarrow Layers \rightarrow Clusters \rightarrow groups

Making a groups or clusters as per the requirements.

eg:-

Gender $\left\{ \begin{array}{l} \rightarrow \text{Male} \\ \rightarrow \text{Female} \end{array} \right.$

Education $\left\{ \begin{array}{l} \rightarrow \text{High school} \\ \rightarrow \text{degree} \\ \rightarrow \text{masters} \\ \rightarrow \text{Phd} \end{array} \right.$

c) Systematic sampling :- Select every n^{th} individual out of Population (N) eg:- every 5th individual selection for credit card selling.

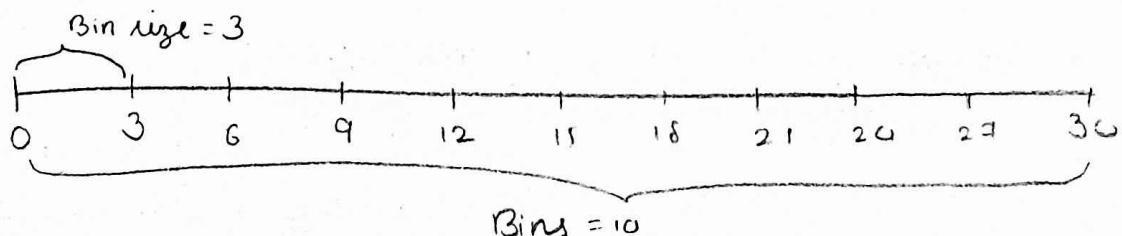
d) Convenience sampling :- Only those who are interested in the survey will only participate.

2) Construction of Histogram?

→ a) Sort the numbers

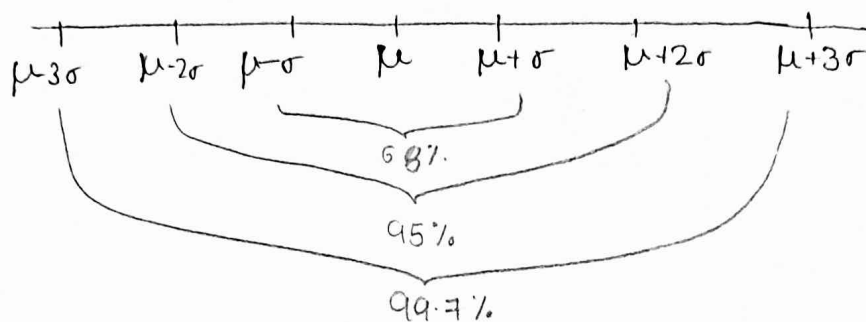
b) Bins \rightarrow No. of groups = 10 (from diagram)

c) Bin size $\rightarrow \frac{\text{Max} - \text{Min}}{\text{Bins}}$ if marks from 0 or $\frac{\text{Max}}{\text{Bins}}$ if not starting from 0 = 3 (from diagram)



3] Empirical Rule of Normal distribution

→



1st standard deviation = 68%.

2nd standard deviation = 95%.

3rd standard deviation = 99.7%.

4] From Gaussian distribution → Standard Normal distribution
 $X \sim \text{Gaussian distribution } (\mu, \sigma)$ → $Y \sim \text{Standard Normal distribution } (\mu=0, \sigma=1)$

$$Z\text{-score} = \frac{x_i - \mu}{\sigma/\sqrt{n}} \quad \text{where } \frac{\sigma}{\sqrt{n}} = \text{standard error}$$

$$= \frac{x_i - \mu}{\sigma} \quad \text{where } n=1$$

$n = \text{sample size}$

5] What is p-value?

→ P-value are used to make a decision about a hypothesis test. P-value is the minimum significant level at which you can reject the null hypothesis. The lower the p-value, the more likely you reject the null hypothesis.

6] How can we relate standard deviation and variance?

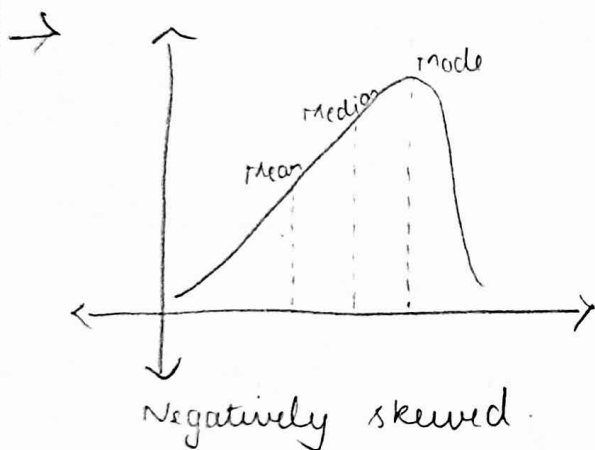
→ Standard deviation refers to the spread of your data from the mean. Variance is the average degree to which each point

differ from the mean i.e the average of all data points
 We can relate standard deviation and variance because it is the square root of variance

7] What is central limit theorem?

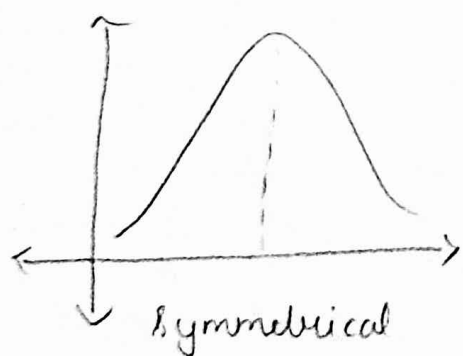
→ The central limit theorem states that if you have a population with mean μ and standard deviation σ and take sufficiently large random samples from the population with replacement, then the distribution of sample means will be approximately normally distributed.

8] Relation between mean, median, mode?

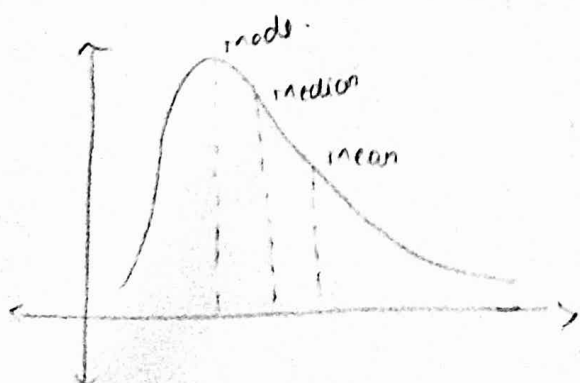


$$\text{mean} < \text{median} < \text{mode}$$

due to large number of skewness in left hand side mean is larger on it



$$\text{mean} = \text{median} = \text{mode}$$

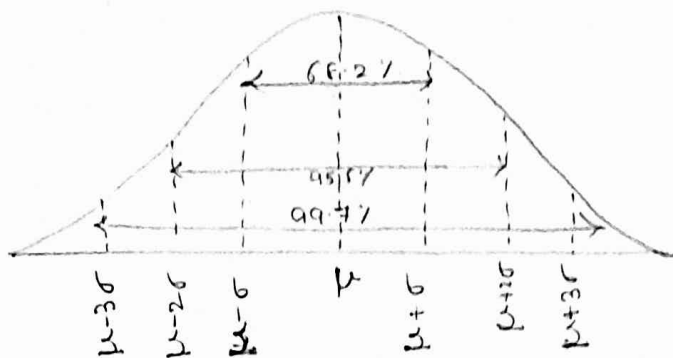


$$\text{mean} > \text{median} > \text{mode}$$

It is slightly skewed data in Right hand side large number of data is found so mean is high on it

q) What is gaussian distribution ?

→



10] What is difference between Standardization and Normalization?

→ In standardization :- $\mu=0, \sigma=1$ by $z\text{-score} = \frac{x_i - \mu}{\sigma}$

In Normalization :- [lower scale, higher scale] is defined by user

$$\text{min max scaler} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

eg [0, 1] → its defined by user.

11] What is KDE Plot?

→ a) A Non-Parametric way of estimation to get Probability density function.

b) Basically used for data smoothing

c) Non-Parametric - Less restriction, Less assumption

Mathematically :-

$$\hat{f}(x, h) = \frac{1}{nh} \sum_{i=1}^n K\{(x - x_i)/h\}$$

$$x_i = \{x_1, x_2, \dots, x_n\}$$

K = kernel function. (Non-Negativity, integrates to 1)

h = smoothing Parameter (bandwidth)

d] Histogram itself a non-parametric But KDE and Histogram is almost similar. But KDE is more interpretable, suggestive, stability to show distribution of data

12] Why sample variance is divided by $\frac{1}{n-1}$?

$$S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

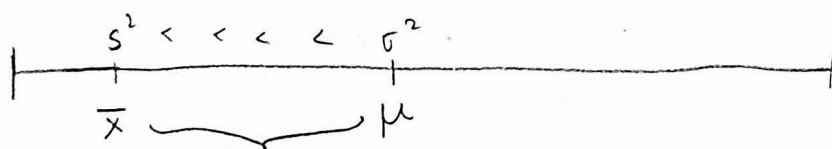
$$\rightarrow \text{a) } S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})}{n}$$

↓
since n is large
 S^2 is small

$$\text{b) } S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})}{n-1}$$

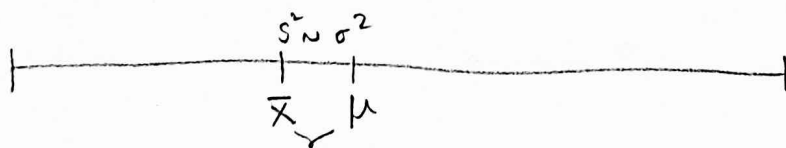
↓
since $n-1$ is small comparatively
 S^2 is large

b) while doing estimation in both the formulas with Population mean, sample mean, Population variance, sample variance.



distance will be large if we apply $\sum_{i=1}^n \frac{(x_i - \bar{x})}{n}$

But



distance will be almost same if we apply $\sum_{i=1}^n \frac{(x_i - \bar{x})}{n-1}$

13) Why to use Z-score?

→ To bring the feature in a ~~small~~ same scale we use Z-score formula

$$\begin{aligned} \text{Z-score} &= \frac{x_i - \mu}{\sigma/\sqrt{n}} \quad \text{where } \frac{\sigma}{\sqrt{n}} = \text{standard error.} \\ &= \frac{x_i - \mu}{\sigma} \quad \text{where } n=1 \end{aligned}$$

14) What is Pearson Correlation Coefficient and Spearman Rank Correlation Coefficient?

→ a) Pearson Correlation coefficient:-

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

i] Strength :- How strongly it is correlated.

ii] direction of Relationship

iii] $-1 \leq \rho \leq 1$

b) Spearman Rank Correlation coefficient?

$$\rho_s = \frac{\text{Cov}[R(X), R(Y)]}{\sigma[R(X)] \times \sigma[R(Y)]}$$

$R(X)$ = Rank of X

$R(Y)$ = Rank of Y

a) Non-linear data.

b) Actually correlation between X and Y.

It's more efficient than Pearson correlation coefficient.

15] What is the difference between a ^{Parameter and} hyperparameter?

→ Model parameters are estimated from data automatically and model hyperparameters are set manually and are used in process to help estimate model parameters.

16] What is the difference between Estimation, Estimator and Estimate

→ Estimation - Process.

Estimator - Formula

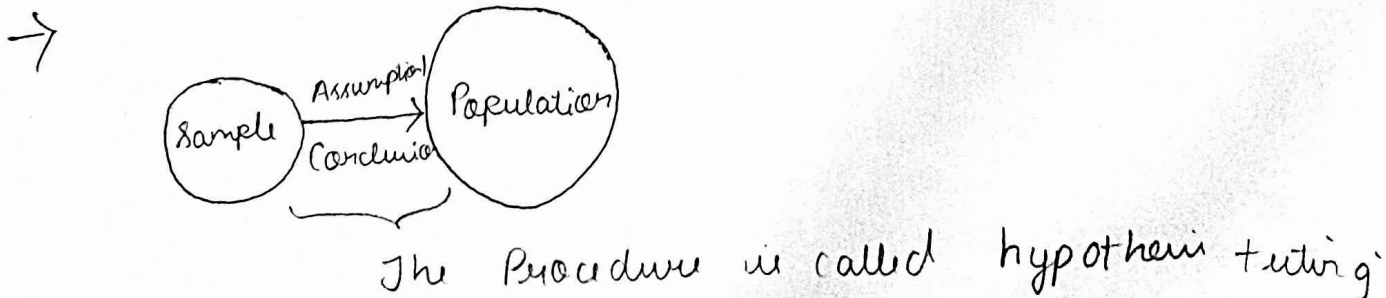
Estimate - Final numerical value.

Estimation :- A Process in which we obtain the values of unknown population parameters with the help of sample data.

Estimator :- It is a rule, formula or function that tells how to calculate an estimate i.e. to estimate the value of a Population parameter.

Estimate - An estimate is the numeric value of the estimator.

17] What is Hypothesis Testing?



18] What is Confidence interval, significance value, p-value?

→ a) Confidence Interval :- It is the probability that a population parameter will fall between a set of values for a certain proportion of times.
for eg C.I = 95%.

b) Significance level \Rightarrow S.V = $1 - C.I = 1 - 0.95 = 0.05$

Confidence interval and significance level are done by the domain expert.

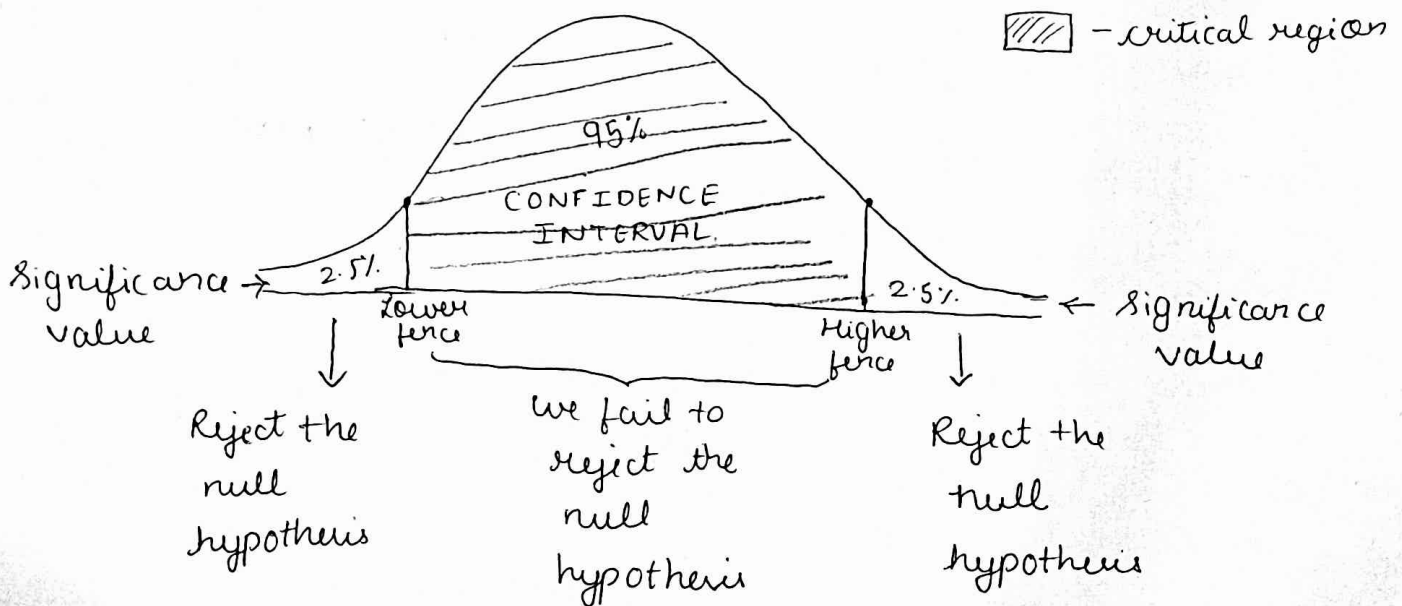
c) P-value :- Tells us how likely it is that your data could have occurred under the null hypothesis.

eg:- $P \leq \alpha$

Reject Null hypothesis

$P > \alpha$

fail to reject null hypothesis.



19] Hypothesis testing Z-test?

→ a] It is used when ($n > 30$)

b] To determine the difference between the Population mean (μ) and sample mean (\bar{x})

c] Population S.D (σ) is known.

d] Used to compare the mean of two samples.

e] One sample : $z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$

$$\text{two sample : } z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

20] Hypothesis testing t-test?

→ a] It is used when ($n < 30$)

b] To determine the difference between sample mean (\bar{x}) and Population mean (μ)

c] Population S.D (σ) is not known.

d] Used to compare the mean of two samples.

e] One sample $t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$

$$\text{two samples } t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$