

Probability

★ An event is a subset of sample space.

★ If $A \cap B = \emptyset$ then $A \cup B$ are mutually exclusive or disjoint sets.

★ The no. of permutations of n distinct objects taken r at a time is ${}^n P_r = \frac{n!}{(n-r)!}$

★ The no. of combinations of n distinct objects taken r at a time is ${}^n C_r = \frac{n!}{r!(n-r)!}$ or denoted by $\binom{n}{r}$.

$$\boxed{P(A \cup B) = P(A) + P(B) - P(A \cap B)} \rightarrow \text{Additive rule.}$$

and if mutually exclusive then $P(A \cup B) = P(A) + P(B)$
bcz $P(A \cap B) = \emptyset$.

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

$$P(A) + P(A') = 1 \quad \text{complement rule}$$

$$P(A \cap B \cap C)$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \quad \text{where } P(A) > 0$$

$$= P(A) \cdot P(B|A)$$

$$P(\bar{A} \cup \bar{B}) = P(\bar{A}) + P(\bar{B}) - P(\bar{A} \cap \bar{B}) \rightarrow P(C|A \cap B)$$

$$P(B|A) = P(B) \quad \text{or} \quad P(A|B) = P(A) \rightarrow \text{independent event}$$

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

$$P(\bar{A} \cap \bar{B}) = 1 - P(A \cup B)$$

$$P(A \cap B) = P(A) P(B) \rightarrow \text{independent event}$$

Bayes tho:- $P(B_r|A) = \frac{P(B_r \cap A)}{\sum_{i=1}^n P(B_i \cap A)}$

$$= \frac{P(B_r) P(A|B_r)}{\sum_{i=1}^n P(B_i) P(A|B_i)}$$

$$\text{for } r = 1, 2, \dots, k$$



$$\text{Mean } \bar{x} = \sum_{i=1}^n \frac{x_i}{n} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

if data is in increased order

$$\bar{x} = \begin{cases} x_{(n+1)/2}, & \text{if } n \text{ is odd,} \\ \frac{1}{2} (x_{n/2} + x_{n/2+1}), & \text{if } n \text{ is even} \end{cases}$$

Sample variance denoted by s^2 ,

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

Sample standard deviation $\Rightarrow s = \sqrt{s^2} = \sqrt{\text{Variance}}$

Median = average of $\left(\frac{n}{2}\right)$ & $\left(\frac{n}{2} + 1\right)$ th value

$$= \begin{cases} \frac{n}{2} \text{ & } \left(\frac{n}{2} + 1\right) & \text{for } n \text{ even} \\ \frac{n+1}{2} & \text{for } n \text{ odd} \end{cases}$$

Trimmed mean :- Suppose 50 data is given and you have to calculate 10% trimmed mean $\rightarrow \frac{10}{100} \times 50 \Rightarrow 5$

\therefore Remove 5 largest values and remove 5 smallest values

Range = max val - min val

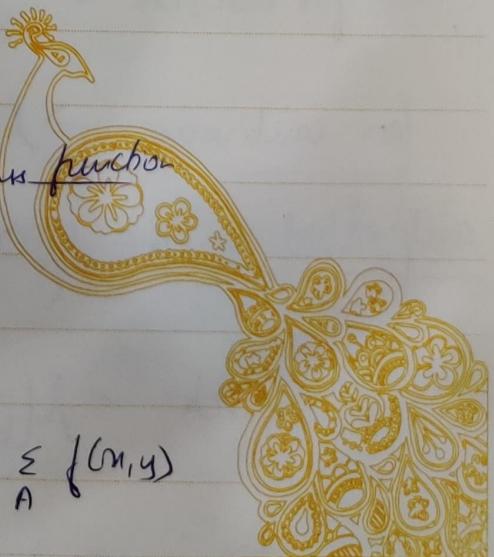
Joint probability distribution or probability mass function

$$\textcircled{1} \quad f(x, y) \geq 0 \text{ for all } (x, y)$$

$$\textcircled{2} \quad \sum_x \sum_y f(x, y) = 1$$

$$\textcircled{3} \quad P(X=x, Y=y) = f(x, y)$$

for any region A in x-y plane $P[(x, y) \in A] = \sum_x \sum_y f(x, y)$



Joint density function

① $f(x,y)$ is jdf of continuous random variables x, y if
 $\Rightarrow f(x,y) \geq 0$ for all (x,y)
 $\Rightarrow \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x,y) dy dx = 1$

$$② P\{(x,y) \in A\} = \iint_A f(x,y) dx dy$$

Marginal distributions of x alone & y alone

$$g(x) = \sum_y f(x,y) \text{ and } h(y) = \sum_x f(x,y)$$

for discrete case and

$$g(x) = \int_0^{\infty} f(x,y) dy \quad \text{and} \quad h(y) = \int_{-\infty}^{\infty} f(x,y) dx \quad (\text{for continuous case})$$

Conditional distribution

$$f(y|x) = \frac{f(x,y)}{g(x)}, g(x) > 0 \quad \text{and} \quad f(x|y) = \frac{f(x,y)}{h(y)}, h(y) > 0$$

$$P(a < x < b | y=y) = \sum_{a < x < b} f(x|y) \quad P(a < x < b | y=y) = \int_a^b f(x|y) dx$$

Statistical independence

Mathematical Expectation

Let x → random variable with probability distribution $f(x)$.

mean or expected value :-

$$\text{for discrete} \rightarrow \mu = E(x) = \sum_n x_n f(x_n) \quad \text{or} \quad \sum_n y_n f(y_n)$$

$$\text{for continuous} \rightarrow \mu = E(x) = \int_{-\infty}^{\infty} x f(x) dx$$

expected value of random variable -

$$\sum_n g(x_n) f(x_n)$$

$$(\text{for discrete}) = M_g(x) = \{g(x)\} = \sum_n g(x_n) f(x_n)$$

$$(\text{for continuous}) = " " = \int_{-\infty}^{\infty} g(x) f(x) dx$$

Correlation coefficient

$$\rho_{xy} = \frac{\text{Cov}(x,y)}{\sqrt{\text{Var}(x)\text{Var}(y)}}$$

Let X, Y be random variables with joint prob dist $f(x, y)$.

The mean expected of random variable $g(x, y)$ is

$$\text{for discrete} \rightarrow M_g(x, y) = E[g(x, y)] = \sum \sum g(x, y) f(x, y)$$

for continuous \Rightarrow

$$" = \int \int g(x, y) f(x, y) dx dy$$

Variance: - for discrete $\Rightarrow \sigma^2 = E[(x - \mu)^2] = \sum (x - \mu)^2 f(x)$ for joint prob

for continuous $\Rightarrow " = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$

\Rightarrow Variance of random variable x is $\sigma^2 = E(x^2) - \mu^2$ Covariance

Discrete Probability Distribution

Binomial distribution :-

Bernoulli trials

$$\text{or } b(x; n, p)$$

(r or x)

The mean and the variance of binomial distribution

$b(x; n, p)$ are $\mu = np$ and $\sigma^2 = npq$ (or $npc(1-p)$)

Multinomial distribution :- ~~X~~

$$f(r_1, r_2, \dots, r_k; p_1, p_2, \dots, p_k, n) = \frac{n!}{r_1! r_2! \dots r_k!} (p_1)^{r_1} (p_2)^{r_2} \dots (p_k)^{r_k}$$

and

$$\sum_{i=1}^k r_i = n \quad \sum_{i=1}^k p_i = 1$$



Hypergeometric Distribution :-

Hypergeometric random variable x , the number of successes in a random sample of size n selected from N items of which k are labelled success and $N-k$ labeled failure, is

$$h(x; N, n, k) = \frac{\binom{k}{x} \binom{N-k}{n-x}}{\binom{N}{n}}, \text{ max } \{0, n - (N-k)\} \leq n \leq \min \{N, k\}$$

The mean and the variance of the hypergeometric distribution $h(x; N, n, k)$ are

$$\mu = \frac{nk}{N} \quad \text{and} \quad \sigma^2 = \frac{N-n}{N-1} \cdot n \cdot \frac{k}{N} \left(1 - \frac{k}{N}\right)$$

$$\therefore \mu = np = \frac{nk}{N} \quad \text{and} \quad \sigma^2 = npq = n \cdot \frac{k}{N} \left(1 - \frac{k}{N}\right)$$

$\frac{N-n}{N-1}$ is negligible

when n is small relative to N .

Multivariate hypergeometric distribution :-

Poisson ratio distribution :-

$\mu = \frac{(\lambda t)^k}{k!} e^{-\lambda t}$ Poisson random variable $x \rightarrow$ no. of outcomes occurring in a given time interval or specific region denoted by t is

$$e^{\mu} \frac{\mu^n}{n!} \rightarrow \mu = \lambda t \quad P(x; \lambda t) = \frac{e^{-\lambda t} (\lambda t)^x}{x!}, \quad x = 0, 1, 2, \dots$$

where λ is the average no. of outcomes per unit time, distance, area or volume and $e = 2.71828 \dots$

Both the mean and variance of poisson distribution $p(x; \lambda)$ is λ .

Continuous probability distribution

Uniform distribution:-

The density function of the continuous uniform random variable X on interval $[A, B]$ is

$$f(x; A, B) = \begin{cases} \frac{1}{B-A}, & A \leq x \leq B \\ 0, & \text{elsewhere} \end{cases}$$

Mean & Variance of uniform distribution

$$\text{is } \mu = \frac{A+B}{2} \text{ and } \sigma^2 = \frac{(B-A)^2}{12}$$

Normal distribution :-

$$n(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

$$\pi = 3.14159 \text{ and } e = 2.71828$$

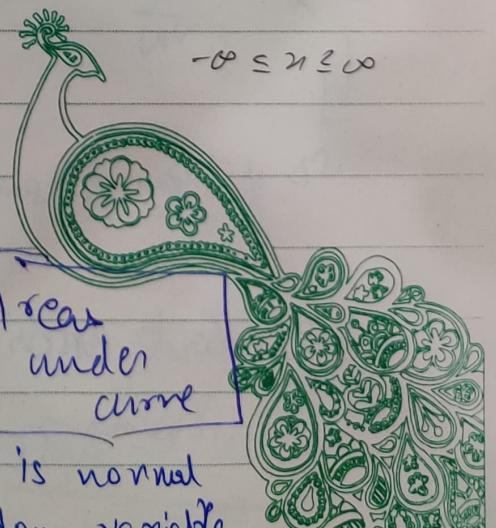
$$\text{Mean} = np$$

$$\text{Standard deviation} = \sigma = \sqrt{npq}$$

$$\therefore Z = \frac{X - \mu}{\sigma}$$

Z is normal random variable with mean 0 & var ≥ 1

A \star area under curve



Normal approximation :-

If X is a binomial random variable with mean $\mu = np$ and variance $\sigma^2 = npq$, then the

limiting form of distribution of

$$Z = \frac{X - np}{\sqrt{npq}}$$

as $n \rightarrow \infty$, is the standard normal distribution

Central limit theorem.

If \bar{X} is the mean of a random sample of size n taken from a population with mean μ and finite variance σ^2 , then the limiting form of the distribution of

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \quad \text{as } n \rightarrow \infty, \text{ is the standard}$$

normal distribution $n(z; 0, 1)$

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} \rightarrow \text{SD of observed pop.} \quad \frac{\bar{X} - \mu_{\bar{X}}}{\sigma_{\bar{X}}}$$

SD of mean for sampling

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

Standard error of
Means:-

$$SEM: - \frac{\sigma}{\sqrt{n}}$$

Interguartile
quartiles

$$\text{Level of confidence} = 100(1-\alpha) = \text{say } 75 \\ \alpha = 0.05$$

range of $|\bar{x} - \mu|$

$$|\bar{x} - \mu| \leq Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

$$\text{Max error } f = Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

Sampling Distribution

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

$$\text{Sample median} := \begin{cases} x_{(n+1)/2} & \rightarrow n \text{ odd} \\ \frac{1}{2} (x_{n/2} + x_{(n/2)+1}) & \rightarrow n \text{ is even} \end{cases}$$

$$\text{Sample variance} := s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$s^2 = \frac{1}{n(n-1)} \left[n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2 \right]$$

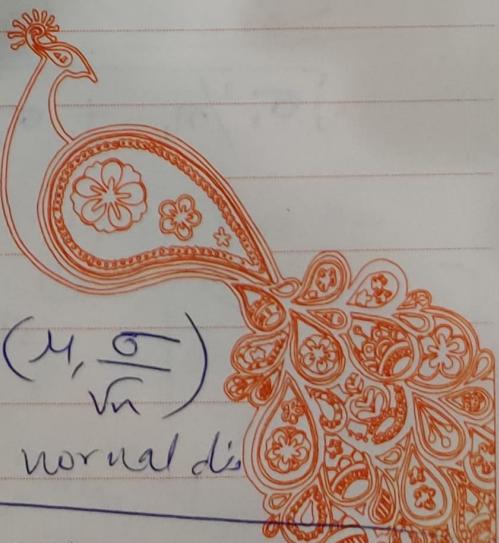
$$\text{Range} = x_{\max} - x_{\min}$$

$$\text{mean} = M \quad \text{var} = \frac{\sigma^2}{n}$$

$$(1) \text{ when } \sigma \text{ known} \rightarrow \bar{x} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

$$(2) \text{ when } \sigma \text{ unknown} \rightarrow \text{pop follows normal dist}$$

t-distribution is after at last



Hypothesis

$H_0 \rightarrow$ null hypothesis $H_1 \rightarrow$ alt. hypothesis \rightarrow claim

Type I error \rightarrow Rejection of H_0 when it's true \rightarrow Prob (Type I error) \leftarrow level of significance

Type II error \rightarrow Non rejection of H_0 when it's false $\rightarrow \beta$

• Prob of Type I error $\propto \frac{1}{\text{Prob of Type II error}}$

• greater n , will reduce α, β

Power:- Prob. of rejecting H_0 given that H_1 is true

$$\downarrow 1 - \beta$$

Level of confidence + level of sign = 1

Test statistics :-

$$Z = \left| \frac{\bar{x} - M}{\frac{\sigma}{\sqrt{n}}} \right|$$

$$\frac{\hat{P} - P_0}{\sqrt{\frac{P_0(1-P_0)}{n}}}$$

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (M_1 - M_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

↓
diff of
two means
 \downarrow
 σ is known

$$S_p^2 = \frac{s_1^2(n-1) + s_2^2(n_2-1)}{n_1+n_2-2}$$

↓
but
 σ is unknown

- ④ know → 2 dust
- ⑤ unknown → +

Conclusion of hypothesis-

Reject H_0 or fail to reject H_0 .

\Rightarrow Level of significance (α)

* Standard value of $\alpha = 5\%$

One Sample

1) Mean (μ)

2) Variance (σ^2)

2 samples

1) Diff in means

2) Ratios of var

More than 2 samples

Multiple means

Anova

for mean (μ):- ① critical Region approach

i) $n < 30 \rightarrow t - \text{dist}^n$

$\text{O}_{ij} \quad n \geq 30 \rightarrow 2\text{-dist}^n$

$\alpha_{1/2}$ or $2\alpha_{1/2}$) or (t_2 or $2t_2$)

where null hypothesis
is rejected

$$(t_{\alpha/2} \text{ or } 2t_{\alpha/2}) \text{ or } (t_{\alpha/2}$$

$(t_{\alpha/2} \text{ or } 2t_{\alpha/2}) \text{ or } (t_{\alpha} \text{ or } 2t_{\alpha})$

→ for two fail typ critical region

$$CR := (-\infty, -t_{\alpha/2}) \cup (t_{\alpha/2}, \infty) ; \quad n < 30$$

$$(-\infty, -2\alpha_{12}) \cup (\alpha_{12}, \infty); \quad n \geq 30$$

\rightarrow for 1 tailed hyp

$\text{CH} \vdash (-\alpha, -\ell_\alpha) \text{ or } (-\alpha, -2\ell_\alpha)$

less than
signs }

(t_0, ∞) or (z_0, ∞) {for greater than}

6

$$\boxed{\begin{array}{l} \text{if } M > M_0, \quad z > z_a \\ M < M_0, \quad z < -z_a \\ M \neq M_0, \quad z < -z_a \quad \text{if } z > z_a/2 \end{array}}$$

P-value approach

if $P\text{ value} < \alpha \Rightarrow$ Reject H_0

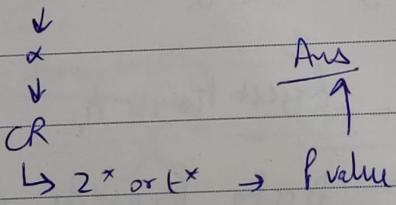
$P\text{ value} \geq \alpha \Rightarrow$ fail to Reject H_0

$$P\text{-value} = P(z < z^*) \xrightarrow{\substack{\downarrow \\ \text{Test stats}}} \text{for left tail test}$$

$$P(z > z^*) \xrightarrow{\substack{\uparrow \\ \text{Test stats}}} \text{for right tail test}$$

$$= 2 [P(z \leq z^*)] \xrightarrow{\substack{\uparrow \\ \text{Test stats}}} \text{for 2 tail}$$

starting of quest H_0, H_1



C II \rightarrow 1 sample $\alpha = \text{significance}$
 $1-\alpha = \text{confidence}$

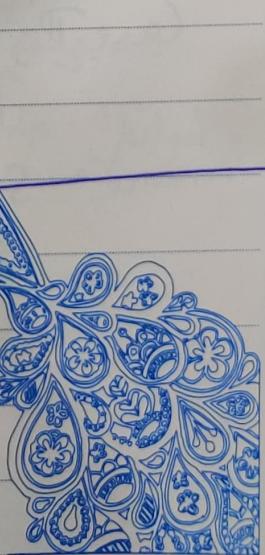
Marginal error $Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

$$(1-\alpha)\% \text{ CI of } \mu = \left(\bar{x} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right)$$

$$n = \left(\frac{Z_{\alpha/2} \cdot \sigma}{M_E} \right)^2$$

$$\hat{\mu}_1 - \hat{\mu}_2 = \bar{x}_1 - \bar{x}_2 - Z_{\alpha/2} \left(\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right),$$

$$|| || + Z_{\alpha/2} \quad || \quad ||$$



Inference of diff of two means

Case I :- when pop SD is known

$$\mu_1, \mu_2, \bar{x}_1, \bar{x}_2, \sigma_1, \sigma_2, n_1, n_2$$

$$t_{cd} = \mu_1 - \mu_2$$

$$H_0 = \mu_0 = \mu_d$$

$$H_1 = \mu_d \neq \mu_0, < \mu_0, > \mu_0$$

$$z^* = \frac{(\bar{x}_1 - \bar{x}_2) - \mu_d}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$



Sample

$$\begin{matrix} \circ \\ \bar{x} \\ \underline{s^2} \end{matrix}$$

$$\mu_1 - \mu_0 < \mu_0$$

$$\mu_1 - \mu_2 > \mu_0$$

$$\mu_1 - \mu_2$$

$$P\text{-value} = 2P(z < z^*) = \alpha < \alpha$$

↓
Reject H_0 in H_1

Ans is statem → Data provide sufficient evidence
to reject H_0 in favor of H_1 .

Case II when SD is unknown

$$\text{pooled Variance} = S_p^2 = \frac{(n_1-1)\bar{s}_1^2 + (n_2-1)\bar{s}_2^2}{n_1+n_2-2} \rightarrow \bar{s}_1 = \bar{s}_2$$

$$\mu_1 - \mu_2 < \mu_0 / t > t_\alpha$$

$$\bar{s}_1^2 = \frac{1}{n_1-1} \sum (\bar{x}_i - \bar{x})^2$$

$$\mu_1 - \mu_2 > \mu_0 / t > t_\alpha$$

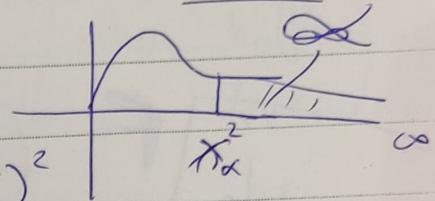
$$\mu_1 - \mu_2 < \mu_0 / t < -t_{\alpha/2}$$

$$T = \frac{(\bar{x}_1 - \bar{x}_2) - \mu_d}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Inference on Pop Variance (or SD)

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

$$V = n-1$$



Chi-Squared dist

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

Single sample

$$H_0: \sigma^2 = \sigma_0^2$$

$$H_1: \sigma^2 \neq \sigma_0^2, < \sigma_0^2, > \sigma_0^2$$

$$\sigma \rightarrow \chi^2_\alpha \text{ or } \chi^2_{1-\alpha}$$

OR

Sample summary
Test stats $\chi^2 = \frac{(n-1)S^2}{\sigma^2}$

$$P \frac{H_0 \text{ value}}{\text{observed value}} \sigma^2 = \sigma_0^2$$

$$\text{upper tailed } H_1: \sigma^2 > \sigma_0^2$$

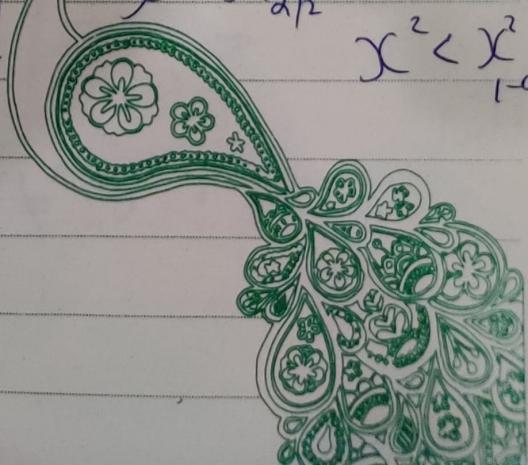
$$\chi^2 > \chi^2_\alpha$$

$$\text{lower tailed } H_1: \sigma^2 < \sigma_0^2$$

$$\chi^2 < \chi^2_{(1-\alpha)}$$

$$\text{two tailed } H_1: \sigma^2 \neq \sigma_0^2$$

$$\chi^2 > \chi^2_{1-\alpha/2} \text{ or } \chi^2 < \chi^2_{\alpha/2}$$



Type 1 - Type 2

Type 1 error

α = (Reject H_0 / when H_0 is true)
 ↓
 CR mai lie ↓
 krega

e.g. ($Z < \sim$ / a or prob)

Type 2 error

β = (fail to reject H_0 / when H_0 is false)
 ↓
 ch mai lie nahi korne
 waala area

e.g. ($Z < \sim$ / a or prob)

$$\text{Power} = 1 - \beta \rightarrow p$$

