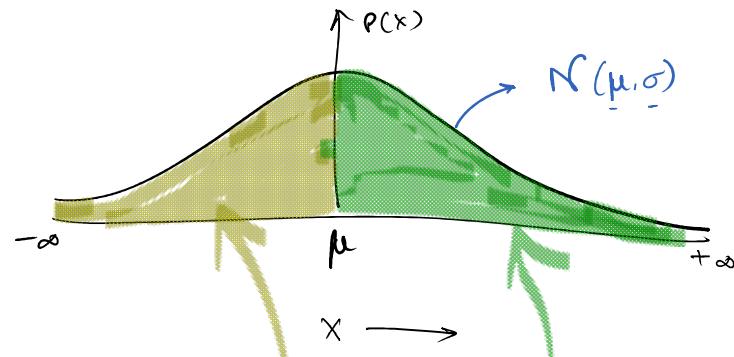


Properties of Normal Distribution



- ① N.D. is a 2-parameter distrib'. \rightarrow mean & std.dev. of the distrib'.
- ② N.D. is symmetric around the mean.

$$\int_{-\infty}^{\mu} P(x) dx = \int_{\mu}^{+\infty} P(x) dx = 0.5$$

Area under the curve on both sides of the mean is same = 0.5.

- ③ mean = median = mode value of N.D.

mean \Rightarrow Normal average

median \Rightarrow Centre point of any distrib' (after sorting)

\rightarrow 50th Percentile \Rightarrow Q₂. (Second Quartile)

mode \Rightarrow Most frequent value of a distrib'. (which value repeats most no. of times).

$$X = \{ 5, 10, 12, 15, 15, 20, 25, 28, 30, 35 \}$$

$$\text{mean} = \frac{1}{10} (5 + 10 + 12 + 15 + 15 + 20 + 25 + 28 + 30 + 35) = \underline{\underline{10}}$$



median $\rightarrow \left(\frac{n+1}{2}\right)^{\text{th}}$ value if n is odd.

→ average of $\frac{n}{2}^{\text{th}}$ & $\left(\frac{n+2}{2}\right)^{\text{th}}$ value. → if n is even

where $n = \text{no. of datapoints}$

$$\frac{10}{2} = 5^{\text{th}} \text{ value}$$

$$\frac{12}{2} = 6^{\text{th}} \text{ value.}$$

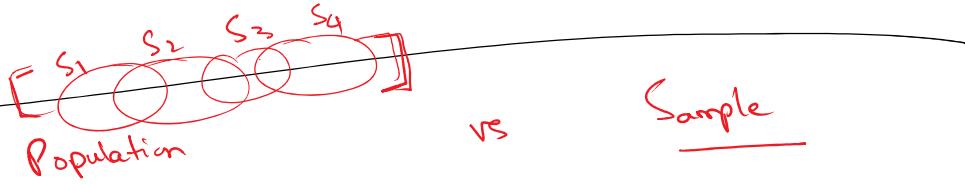
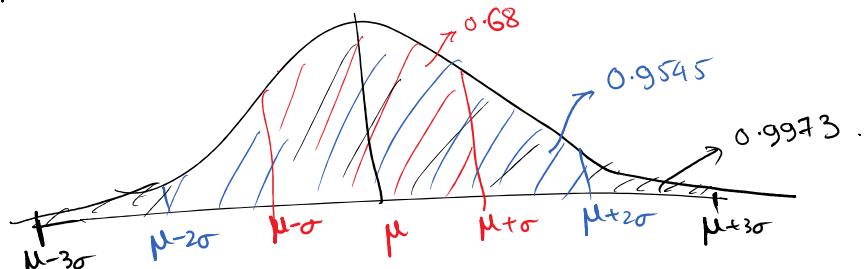
Hence the average of 5th & 6th value = $\frac{15+20}{2} = 17.5$

④ VV Imp:

$$\mu - \sigma \leq x \leq \mu + \sigma \Rightarrow \text{Area under the ND curve} = 0.6828$$

$$\mu - 2\sigma \leq x \leq \mu + 2\sigma \Rightarrow \dots = 0.9545$$

$$\mu - 3\sigma \leq x \leq \mu + 3\sigma \Rightarrow \dots = 0.9973$$



(μ, σ, σ^2) are called as

Population

Parameters

(\bar{x}, s, s^2) are called as
sample characteristics.

mean = \bar{x}

std. dev = s

$$s = \sqrt{\frac{1}{n-1} \cdot \sum (x_i - \bar{x})^2}$$

n = sample size

$$\text{Variance} = s^2 = \frac{1}{n-1} \cdot \sum (x_i - \bar{x})^2$$

Bessel Correction

$$df = n-k$$

$$\text{mean} = \mu$$

$$\text{std. dev} = \sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

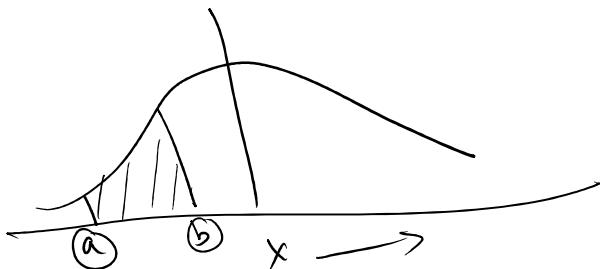
N = Popl. Size

$$\text{Variance} = \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

$$\text{as } n \rightarrow N, s \rightarrow \sigma$$

$$as \quad n \rightarrow \infty, \quad \text{dof} = nk$$

ddof = diff. between n & dof
 $= n - (n-k) = k.$



Area under the curve between a & b

$$= \int_a^b p(x) dx$$

Sum of Prob. of all the values that x can take between a & b .

$$P(a \leq x \leq b) = \int_a^b p(x) dx$$

Cumulative Distribution Function (CDF)

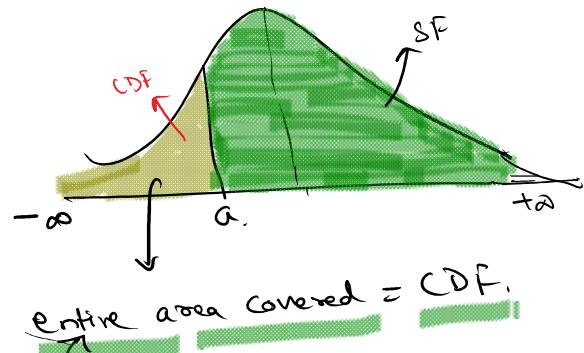
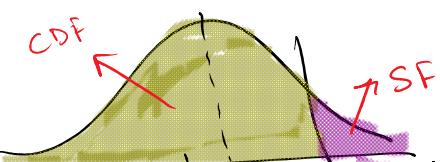
represents the Prob. that a random variable X takes a value less than or equal to a .

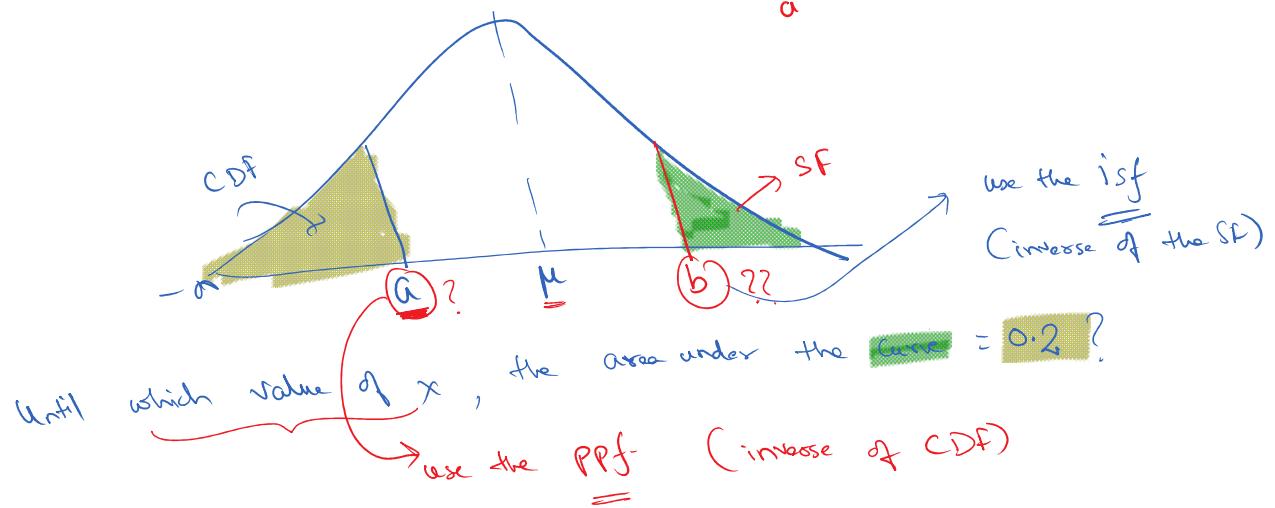
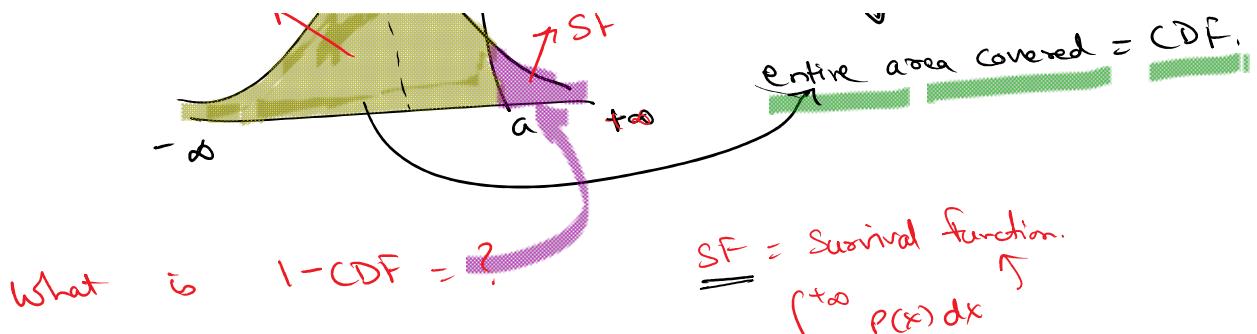
Discrete RV

$$P(X \leq a) = \sum_{k=0}^a P(X=k) = P(X=0) + P(X=1) + P(X=2) + \dots + P(X=a).$$

Cont. RV

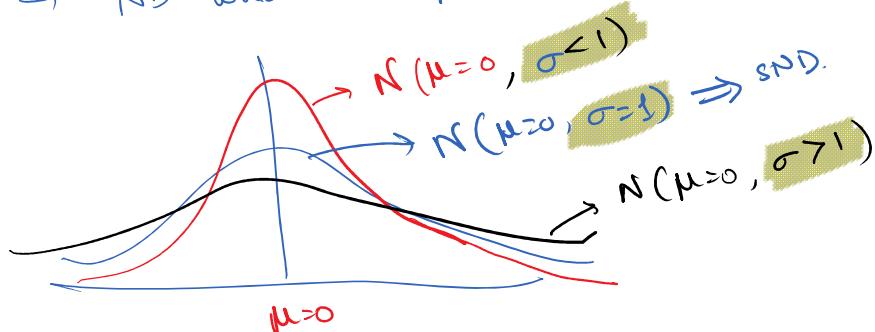
$$P(X \leq a) = \int_{-\infty}^a p(x) dx$$



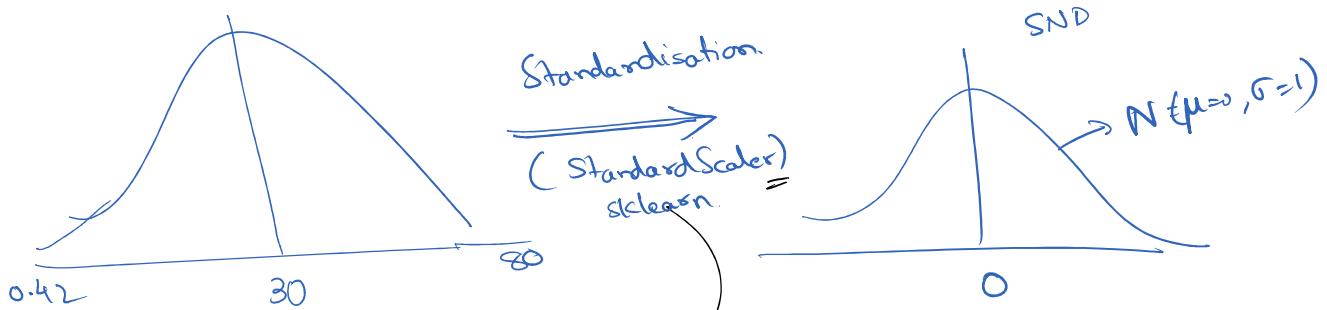


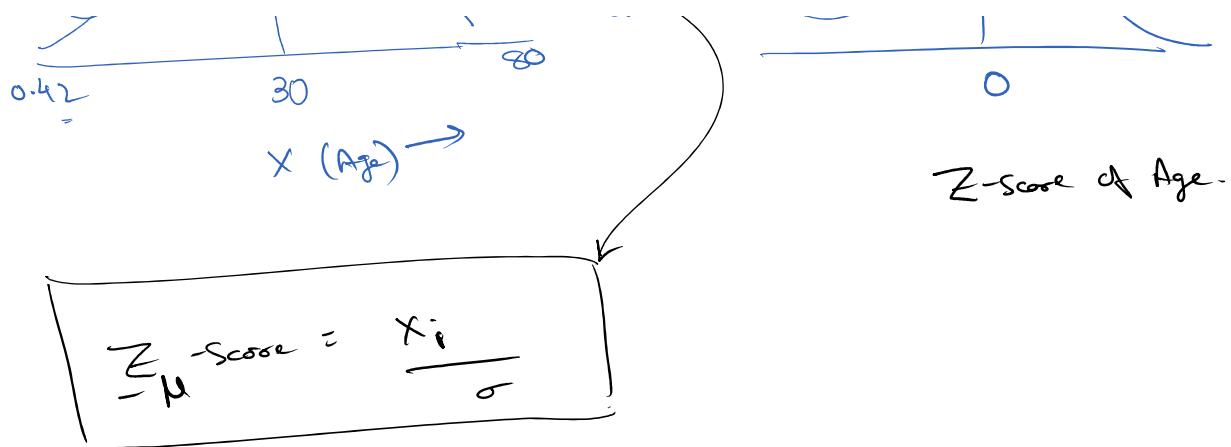
Standard Normal Distribution.

SND. \rightarrow ND which has $\mu=0$, $\sigma=1$.



As the $\sigma \uparrow$, the distribution becomes more flatter / Broader.
 As the $\sigma \downarrow$, the distribution becomes more narrow.

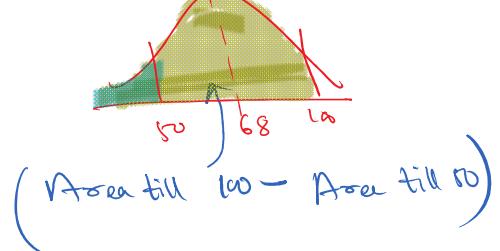
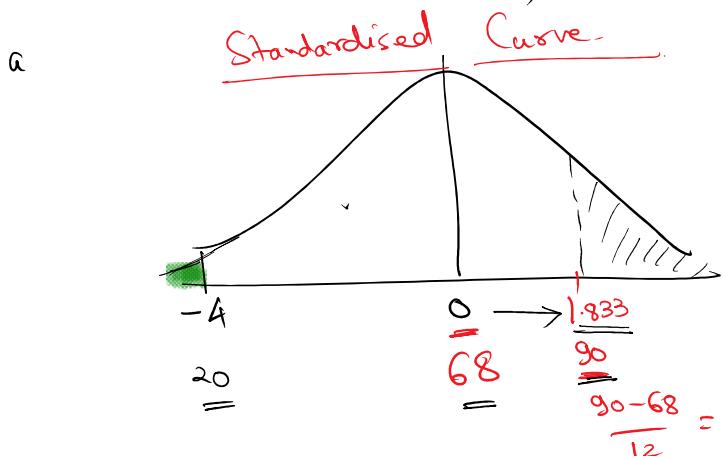
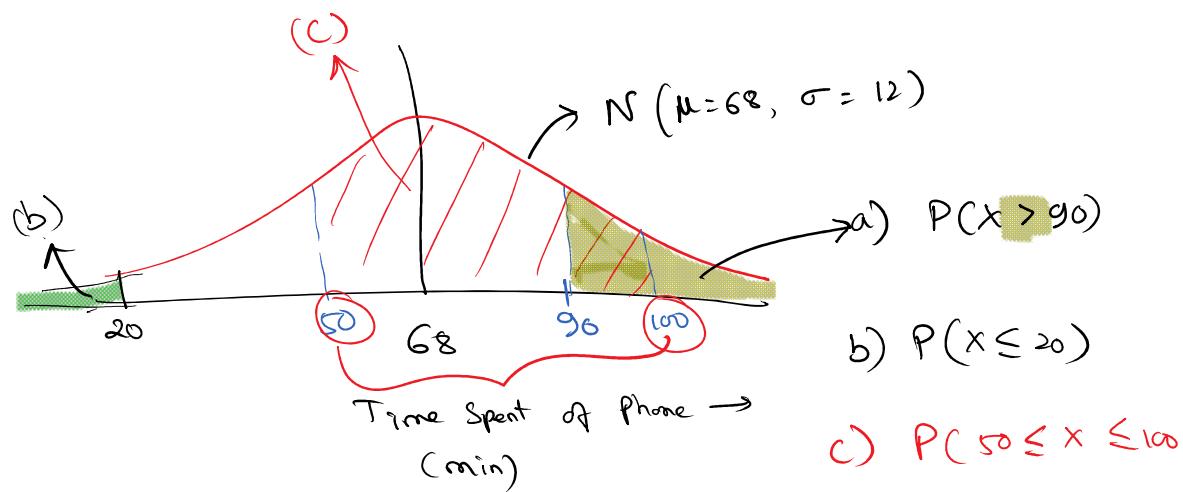




Exercise:

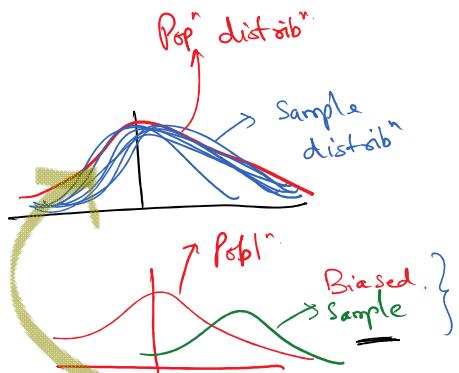
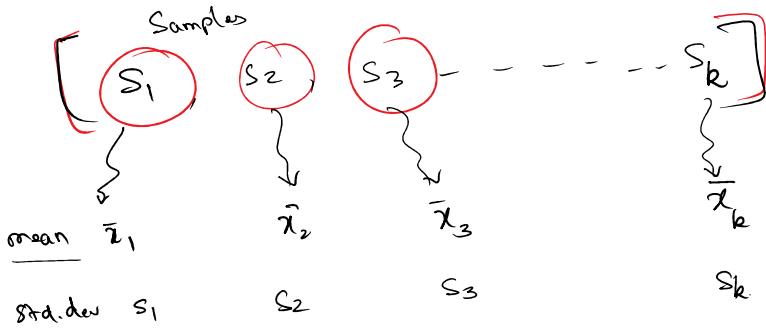
According to a survey on use of smart phones in India, the smart phone users spend 68 minutes in a day on average in sending messages and the corresponding standard deviation is 12 minutes. Assume that the time spent in sending messages follows a normal distribution.

- What proportion of the smart phone users are spending more than 90 minutes in sending messages daily?
- What proportion of customers are spending less than 20 minutes?
- What proportion of customers are spending between 50 minutes and 100 minutes?



$$\frac{1}{12} = 1.855$$

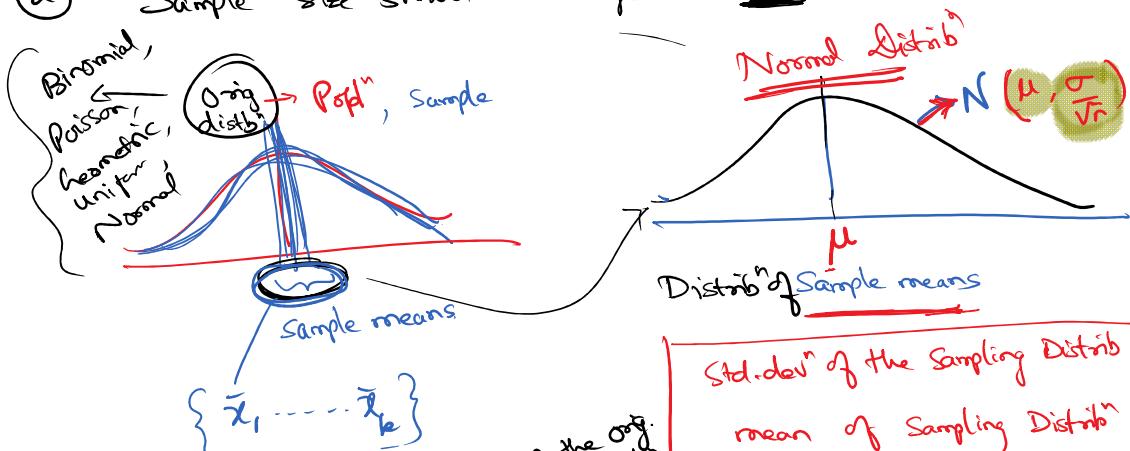
Central Limit Theorem



Assumptions: ① These are IID Samples (Independent & Identically Distributed)

basically means that each sample has been drawn independent of each other \rightarrow Sampling with replacement \Rightarrow Bootstrap Sampling.

② Sample size should be large ($N > 30$)



$$\text{Std.dev of the Sampling Distrib} = \frac{\sigma}{\sqrt{n}}$$

$$\text{mean of Sampling Distrib} = \mu$$

According to CLT: The distib' of sampling means follows Normal Distib' Irrespective of the popl distib'

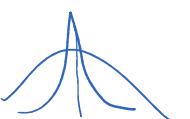
① with its mean equal to the popl mean.

$$\bar{x}_1 + \bar{x}_2 + \bar{x}_3 + \dots + \bar{x}_k = \underline{\underline{\mu}}$$

② its std. dev equal to $\left(\frac{\sigma}{\sqrt{n}}\right)$

$$\left(\frac{\sigma}{\sqrt{n}}\right)$$

n = Sampling size



Nn

Standard Error in the estimation of the Popl' mean form
 the sample means = $\frac{\sigma}{\sqrt{n}}$

$s_1, s_2, s_3, \dots, s_k$

a) We can estimate the popl' parameters. \Rightarrow Confidence Interval Estimation

e.g. $s_1 \rightarrow \bar{x}, s_1$
 popl' mean lies somewhere $\left[\bar{x} \pm Z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right) \right]$

$\left(\bar{x} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right) \quad \bar{x} \quad \left(\bar{x} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right)$

b) You can Hypothesize a value of popl' mean, from the Sample mean & try to validate this number.

e.g. $\bar{x} = 100$, $\left[\mu_0 > 110, \mu_0 < 95, \mu_0 = 110 \right]$
 RT = LT = $2T =$

HYPOTHESIS TESTING

Step 1 : Define H_0 & H_A .

H_0 \rightarrow Null Hypothesis \rightarrow formulated in a way to "Nullify" your H_A .
 H_A \rightarrow Alternate Hypothesis. \rightarrow Whatever you claim against the "Norm" or common understanding.
 (equity sign $\rightarrow H_0$)

Step 2 : Determine which test? Z-test vs t-test, χ^2 -test etc...

If $n < 30$, t-test.

If $n < 30$, **t-test**.
 If std.dev of popl" (σ) is unknown \rightarrow **t-test**.
 Otherwise go for **Z-test**.
 H.T. for **popl" means**.

Step 3 : Cal the Test Statistic..

$$Z_{\text{stat}} = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

μ_0 = Hypothesized popl" mean

$$t_{\text{stat}} = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

$(1-\alpha) \Rightarrow$ confidence Level.

Step 4 : Cal the Critical value. Z_{crit} , t_{crit} . (based on α)

$\alpha \rightarrow$ Significance level.

"**min**" evidence needed in favor of H_0

Step 5 : Decision Making.

Compare Z_{stat} with Z_{crit}
And conclude depending

Left-Tailed.

Right-Tailed.

2-Tailed.

Cal p -value
of compare with α .

p -value (prob/evidence in favor of H_0)

Hence, if $p > \alpha$, "Fail to Reject H_0 "

$p < \alpha$, \Rightarrow Reject H_0 .

actual evidence presented by the sample

$p < \alpha \Rightarrow$ cumulative prob. (sf, cdf)



You (Judge)

Prosecution

(H_A) Person is

Accused of murder

Defense.

H_0 : Person is innocent.

Prosecution

(H_A) Person is murderer

- ① Judge pronounces guilty → Person is murderer
- ✗ ② Judge pronounces guilty → Person is innocent → β
- ✗ ③ Judge sets him free → Person is murderer → α
- ✓ ④ Judge sets him free → Person is innocent

Type-I and Type-II Errors

Type I Error: Conditional probability of rejecting a null hypothesis when it is true is called Type I Error or False Positive (falsely believing that the claim made in alternative hypothesis is true). The significance value α is the value of Type I error.

$$\text{Type I Error} = \alpha = P(\text{Rejecting null hypothesis} \mid H_0 \text{ is true})$$

Type II Error: Conditional probability of failing to reject a null hypothesis (or retaining a null hypothesis) when the alternative hypothesis is true is called Type II Error or False Negative (falsely believing that there is no relationship).

$$\text{Type II Error} = \beta = P(\text{Retain null hypothesis} \mid H_0 \text{ is false})$$

$$\text{Power of the test} = 1 - \beta = 1 - P(\text{Retain null hypothesis} \mid H_0 \text{ is false})$$

$$\text{Alternatively the power of test} = 1 - \beta = P(\text{Reject null hypothesis} \mid H_0 \text{ is false})$$

Condition for rejection of null hypothesis H_0

Type of Test	Condition	Decision
Left-tailed test	$Z\text{-statistic} < \text{Critical value}$	Reject H_0
	$Z\text{-statistic} \geq \text{Critical value}$	Retain H_0
Right-tailed test	$Z\text{-statistic} > \text{Critical value}$	Reject H_0
	$Z\text{-statistic} \leq \text{Critical value}$	Retain H_0
Two-tailed test	$ Z\text{-statistic} > \text{Critical Value} $	Reject H_0
	$ Z\text{-statistic} \leq \text{Critical Value} $	Retain H_0