

□ Continuous Distribution:-

1. Normal Distribution :-

It is also called as Gaussian Distribution.

(i) Definition :-

A random variable x is said to be follow a normal distribution with parameters μ and σ , if its pdf is given by :-

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

for $-\infty < x < \infty$ and
 $-\infty < \mu < \infty$ and
 $\sigma > 0$.

(ii) We also write : $x \sim N(\mu, \sigma^2)$.

(iii) $E(x) = \mu$, $\text{Var}(x) = \sigma^2$.

(iv) If $\mu=0$, $\sigma^2=1$, the x is said to follow a standard Normal Distribution ; $x \sim N(0, 1)$.

(v) The pdf of standard Normal Distribution is given by :

$$\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} ; -\infty < x < \infty.$$

(vi) a) The density of a Normal distribution has it's maximum at $x=\mu$.

b) The density is symmetric around μ .

inflection points of density are

(vii) $\alpha = (\mu - \sigma)$ and $\alpha = (\mu + \sigma)$.

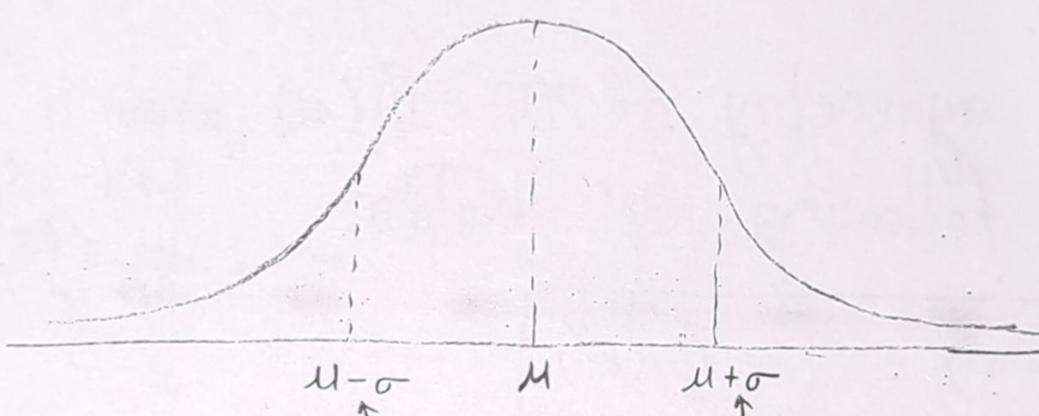
a) A lower value of σ indicates higher concentration around μ .

b) A high value of σ indicates flatter density.

(viii) The CDF of $X \sim N(\mu, \sigma^2)$ is

$$F(x) = \int_{-\infty}^x \phi(t) dt = \Phi(x).$$

(ix) Graph of PDF:-



(x) Calculation rules for Normal Random Variable:

a) Let, $X \sim N(\mu, \sigma^2)$

using the transformation,

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1).$$

Every Normally distributed Random Variable can be transformed into a standard Normal R.V.

We call this transformation as Z transformation.

$$\text{b) } P(X \leq b) = P\left(\frac{X-\mu}{\sigma} \leq \frac{b-\mu}{\sigma}\right)$$

$$= P\left(Z \leq \frac{b-\mu}{\sigma}\right).$$

$$= \Phi\left(\frac{b-\mu}{\sigma}\right).$$

$$\text{c) } P(X > a) = 1 - P(X \leq a)$$

$$= 1 - \Phi\left(\frac{a-\mu}{\sigma}\right)$$

$$\text{d) } P(a \leq X \leq b) = \Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right)$$

e) Because of symmetry of PDF $\phi(x)$ around its mean, the following eqn holds.

$$\boxed{\Phi(-a) = 1 - \Phi(a)}.$$

$$\begin{aligned} P(-a < z < a) &= \\ &P(z \leq a) - P(z \leq -a) \\ &= \Phi(a) - \Phi(-a) \\ &= \Phi(a) + \Phi(a) - 1 \end{aligned}$$

$$\boxed{P(-a < z < a) = 2\Phi(a) - 1.}$$

$$\boxed{= 2\Phi(a) - 1.}$$

R-Commands:-

$$P(X \leq r) = P(X < r) = \underline{\underline{\text{pnorm}(r, \mu, \sigma)}}.$$

For standard normal = pnorm(r, 0, 1).

Q. Let, x denotes the no. of scores in a test. If x is normally distributed with mean 100 and st. deviation 15, find the prob. that x does not exceed 130.

$$\rightarrow \mu = 100, \sigma = 15, Z = \frac{x-\mu}{\sigma}$$

$$x = 130, \therefore Z = 2.$$

$$\therefore P(x \leq 130) = \underline{0.9772} \quad \begin{matrix} \text{From Table, OR} \\ \int_{-\infty}^{130} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx \end{matrix}$$

Q. The random variable x is normally distributed with $\mu = 9$, and $\sigma = 3$. Find probability when

$$(i) x \geq 15.$$

$$(ii) x \leq 15.$$

$$(iii) 0 < x < 9$$

$$(iv) x < 15.$$

$$\rightarrow P(x \geq 15) = P(z \geq 2) = \underline{0.0228}.$$

$$Z = \frac{15-9}{3} = 2$$

$$\therefore P(x \leq 15) = \underline{0.9772} = 1 - P(x \geq 15)$$

$$\therefore P(0 < x \leq 9) = P(x \leq 9) - P(x \leq 0)$$

$$= P(z \leq 0) - P(z \leq -3)$$

$$= 0.5 - 0.0013$$

$$= \underline{0.4987}$$

$$P(x < 15) = \underline{0.9772}$$

Q. A research scientist reports that a mice will live an average of 40 months, when their diets are sharply enriched with proteins and fats. Assuming that the life time of such mice are normally distributed, with $\sigma = 6.3$ months. Find prob. that the given mice will live
 (i) more than 32 months.
 (ii) less than 28 months.
 (iii) more than 37 months and/or less than 49 months.

$$\rightarrow \mu = 40, \sigma = 6.3.$$

$$(i) X = 32. Z = \frac{X-\mu}{\sigma} = \frac{32-40}{6.3} = -1.27.$$

$$\therefore P(X > 32) = 1 - P(X \leq 32) \\ = 1 - P(Z \leq -1.27) = 1 - 0.1020 \\ = \underline{\underline{0.8980}}$$

$$(ii) P(X < 28) = P(Z < -1.90) \\ = 0.0287$$

$$(iii) (37 < X < 49) = P(-0.47 < Z < 1.43) = 0.6044.$$

Q. The mean height of 500 students 151 cm. $\sigma = 15$ cm. Assuming that the heights are normally distributed, find how many student's height lies bet ~~120~~ 120 and 155 cm.

$$\rightarrow \mu = 151, \sigma = 15.$$

$$P(120 < X < 155) = P(-2.06 < Z < 0.27) \\ = \underline{\underline{0.5867}}$$

$$\therefore \text{No. of students} = 0.5867 \times 500 \\ = \underline{\underline{293}}$$

Q. If x is Normal, with $\mu=100$ and $\sigma=5$.

find: a) $P(95 < x < 110)$

b) $P(x < 50)$

c) K , if $P(x > K) = 0.3192$

d) x_1, x_2 if P .

$$\rightarrow \text{a) } P(95 < x < 110) = P(-1 < z < 2)$$

$$= \underline{\underline{0.8185}}$$

$$\text{b) } P(x < 50) = P(z < 10)$$

$$= \underline{\underline{0}}.$$

$$\text{c) } P(x > K) = 0.3192.$$

$$\therefore 1 - P(x \leq K) = 0.6808$$

$$\therefore z = 0.47 = \frac{x - \mu}{\sigma} \quad \therefore \underline{\underline{x = 102.35}}$$

Exponential Distribution:

It is used when we assume that the future lifetime is independent of the lifetime that has already taken place.

Then the waiting time can be considered to be exponentially Distributed.

(i) Defⁿ: - A R.V. X is said to follow an exponential distribution with parameters $\lambda > 0$ if its pdf is given by,

$$f(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{if } x \geq 0. \\ 0, & \text{elsewhere.} \end{cases}$$

Also, $X \sim \text{Exp}(\lambda)$.

(ii) $E(X) = \underline{\frac{1}{\lambda}}$ $\text{Var}(X) = \underline{\frac{1}{\lambda^2}}$

(iii) CDF is,

$$F(x) = P(X \leq a) = P(X < a) = \begin{cases} \underline{1 - e^{-\lambda a}}, & a \geq 0 \\ 0, & \text{else.} \end{cases}$$

Theorem: - The no. of events Y , occurring within a Continuum of time is Poission Distributed, with Parameter λ , IFF, the time betⁿ two events is exponentially distributed with parameter λ .

R-Command:

$$\underline{P(X \leq x) = pexp(x, \lambda)}.$$

Memorylessness Property:-

If time t has already been reached, then the probability of reaching a time greater than $t + \Delta$ does not depend on t :

i.e.
$$\boxed{P(X > t + \Delta | X > t) = P(X > \Delta)}.$$

$$\begin{aligned}
 P(X > t + \Delta | X > t) &= \frac{P(X > t + \Delta \cap X > t)}{P(X > t)} \\
 &= \frac{P(X > t + \Delta)}{P(X > t)} \quad (\Delta > 0) \\
 &= \frac{1 - P(X \leq t + \Delta)}{1 - P(X \leq t)} \\
 &= \frac{1 - (1 - e^{-\lambda(t + \Delta)})}{1 - (1 - e^{-\lambda t})} \\
 &= e^{-\lambda \Delta} = 1 - (1 - e^{-\lambda \Delta}) \\
 &= 1 - P(X \leq \Delta) \\
 \underline{P(X > t + \Delta | X > t)} &= \underline{P(X > \Delta)}.
 \end{aligned}$$

Q. The length of telephone conversation is an exponential variate with mean 3; Find the prob., that a call

(i) Ends in less than 3 min.

(ii) Takes bet'n 3 to 5 min.

$$\lambda = \frac{1}{3}$$

$$\rightarrow (i) P(x \leq 3) = 1 - e^{-\lambda x} = 1 - e^{-1} = \frac{e-1}{e}$$

$$\begin{aligned} (ii) P(3 \leq x \leq 5) &= (e^{-5\lambda}) - (e^{-3\lambda}) \\ &= e^{-1} - e^{-5/3} = \frac{1}{e} \left(1 - \frac{1}{e^{2/3}}\right) \end{aligned}$$

Distribution of Arithmetic Mean for Normal Distribution :-

Assume that $X \sim N(\mu, \sigma^2)$.

Consider a random sample $X = (x_1, x_2, x_3, \dots, x_n)$ of n i.i.d. Random Variable with $x_i \sim N(\mu, \sigma^2)$,
(Independently Identically Distributed),

then \bar{X} = Arithmetic Mean, then

$$\underline{\underline{E(\bar{X}) = \mu}} \quad \text{and} \quad \underline{\underline{\text{Var}(\bar{X}) = \frac{\sigma^2}{n}}} \quad - (\text{Proved Earlier}).$$

$$\underline{\underline{\text{Std. Dev. } (\bar{X}) = \frac{\sigma}{\sqrt{n}}}}$$

Note:- If x_1, x_2, \dots, x_n are n independent normal random variables with mean $\mu_1, \mu_2, \dots, \mu_n$ and variance $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$ then for any real no. a_1, a_2, \dots, a_n

$$\underline{a_1x_1 + a_2x_2 + \dots + a_nx_n \sim N(a_1\mu_1 + a_2\mu_2 + \dots + a_n\mu_n, a_1^2\sigma_1^2 + a_2^2\sigma_2^2 + \dots + a_n^2\sigma_n^2)}$$

holds.

Central Limit Theorem:-

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

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \quad \text{as } n \rightarrow \infty \text{ is Standard Normal Distribution } N(0, 1)$$

- (i) The sample size $n=30$, is a guideline to use for the central limit thm.
- (ii) The normal approximation for \bar{x} will be good if $n \geq 30$, provided that the population distribution is not skewed.
- (iii) The population is not too different from a normal distribution then the approximation is good for $n < 30$ as well.

Sampling Distribution:-

There are 3 types of Sampling distribution.

1. χ^2 -Distribution.
2. t-Distribution.
3. F-Distribution.

χ^2 -Distribution:-

(i) Defⁿ :- Let, z_1, z_2, \dots, z_n be i.i.d. Normal Random Variable. Then, the sum of their squares is

$$z_1^2 + z_2^2 + \dots + z_n^2 = \sum_{i=1}^n z_i^2 \text{ is } \underline{\chi^2\text{-Distributed with}} \\ \underline{\text{degree of freedom } n}. \quad (\text{df} = \text{degree of freedom})$$

$$(ii) f(x) = \begin{cases} \frac{x^{\frac{n}{2}-1} e^{-x/2}}{\Gamma(\frac{n}{2}) 2^{n/2}}, & x \geq 0 \\ 0 & \text{elsewhere} \end{cases}$$

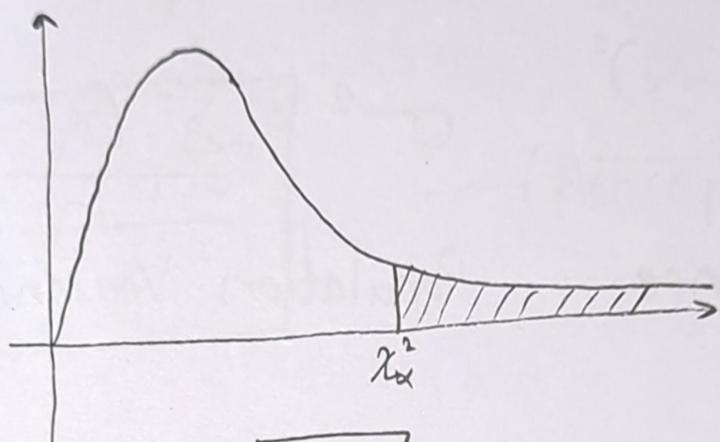
where Γ is gamma function

i.e.

$$\Gamma = \begin{cases} (r-1)! & , \text{ for } r \text{ integer} \\ \int_0^\infty t^{r-1} e^{-t} dt & , \text{ elsewhere.} \end{cases}$$

- (iii) χ^2 -Distribution is not symmetric. (skewed graph).
- (iv) χ^2 -lies betⁿ 0 to ∞ .
- (v) Mean = $E(X) = df$
 $Var(X) = 2 \cdot df$.
- (vi) Small Value of df gives high value of Tail.

Graph:



$$\boxed{\chi^2_{df}}$$

(vii) Theorem :-

Consider two independent random variables which are χ_m^2 and χ_n^2 distributed, then the sum of these two R.V. is χ_{m+n}^2 distributed.

Q. Using stastcial tables find

(i) $\chi_{0.025}^2$ when $df = 15$. $\underset{2}{\sim} 27.488$

(ii) $\chi_{0.05}^2$ when $df = 25$. $\underset{2}{\sim} 37.652$

Q. Find the prob. that a random sample of 25 observations from a normal population with Var. 6 will have a variance (i) s^2 greater than 9.1
(ii) s^2 lies betⁿ 3.462 and 10.745.

P.T.O. \longrightarrow

Theorem :- A χ^2 -distributed random variable, having sample variance S^2 in an i.i.d. sample of size n from a normally distributed population is given by

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$$

Formulas :-

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

Sample Variance.

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

Population Variance.

Proof :-

$$\sum_{i=1}^n z_i = \sum_{i=1}^n \left(\frac{x_i - \mu}{\sigma} \right)^2 = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

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

$$= \frac{1}{\sigma^2} \sum_{i=1}^n \left[(x_i - \bar{x})^2 + (\bar{x} - \mu)^2 + 2(x_i - \bar{x})(\bar{x} - \mu) \right]$$

$$= \frac{1}{\sigma^2} \left\{ \sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n (\bar{x} - \mu)^2 + 2(\bar{x} - \mu) \sum_{i=1}^n (x_i - \bar{x}) \right\}$$

$$= \frac{1}{\sigma^2} \left\{ \sum_{i=1}^n (x_i - \bar{x})^2 + (\bar{x} - \mu)^2 \times n + 2(\bar{x} - \mu)(0) \right\}$$

$$= \frac{1}{\sigma^2} \left\{ \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2 \right\}$$

$$= \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2 + \left(\frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \right)^2$$

$$\sum_{i=1}^n z_i^2 = \frac{(n-1)s^2}{\sigma^2} + \left(\frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \right)^2$$

$$= \frac{(n-1)s^2}{\sigma^2} + Z^2 \quad \text{— (By Central Limit Thm.)}$$

$$\therefore \sum_{i=1}^n z_i^2 - Z^2 = \frac{(n-1)s^2}{\sigma^2}$$

$$\Rightarrow \boxed{\chi_{n-1}^2 \sim \frac{(n-1)s^2}{\sigma^2}} \quad \text{— (Hence, Proved).}$$

Q. Continued. $n=25$. $\sigma^2=6$.

$$\text{Q. } P(s^2 > 9.1) = P \left(\frac{(n-1)s^2}{\sigma^2} > \frac{9.1 \times (n-1)}{\sigma^2} \right)$$

$$= P \left(\chi_{n-1}^2 > \frac{9.1 \times 24}{6} \right) = P \left(\chi_{n-1}^2 > 36.4 \right)$$

$$= \underline{0.05.}$$

$$\text{5) } P(3.462 < s^2 < 10.745) = P(13.848 < \chi_{24}^2 < 42.980)$$

$$= P(\chi_{24}^2 > 13.848) - P(\chi_{24}^2 > 42.980)$$

$$= 0.95 - 0.01$$

$$= \underline{\underline{0.94.}}$$

Q. An electrical firm manufactures light bulbs that have a length of light that is approximately normally distributed with mean = 800 hrs; and st. dev. = 40 hrs. Find the prob. that a random sample of 16 bulbs will have an average life of less than 775 hrs.

$$\rightarrow \mu = 800, \sigma = 40. \quad P(\bar{X} < 775) = ?$$

$$P(\bar{X} < 775) = P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < \frac{775 - 800}{40/4}\right) = P(Z < -2.5)$$

$$= \underline{\underline{0.0062.}}$$

t-Distribution:

(i) Let, X and Y be two independent random variables such that $X \sim N(0, 1)$ and Y is χ_n^2 -distributed, then the ratio $\frac{X}{\sqrt{Y/n}} \sim t_n$.

(ii) pdf is given by :-

$$f(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{\frac{n}{2}} \cdot \sqrt{n\pi}} \left(1 + \frac{x^2}{n}\right)^{-\left(\frac{n+1}{2}\right)}, \quad -\infty < x < \infty.$$

Student's Theorem:-

Let, $X = (X_1, X_2, \dots, X_n)$ with $X_i \stackrel{i.i.d}{\sim} N(\mu, \sigma^2)$
 then the ratio $\left(\frac{\bar{X}-\mu}{S}\right)\sqrt{n} = \frac{(\bar{X}-\mu)\sqrt{n}}{\sqrt{\frac{1}{(n-1)} \sum_{i=1}^n (X_i - \bar{X})^2}}$ $\sim t_n$, i.e.
 is t-distributed with $(n-1)$ degree of freedom.

(iv) If Sample Size, $n < 30$ then we use t-Distribution
 and for $n \geq 30$, we use Normal Distribution.

$$(v) t_\alpha = -t_{1-\alpha}$$

Q. Find t-value with dof 14 leaving an area of 0.025 to the left.

$$\Rightarrow \mu = 14, \quad \alpha = 1 - 0.025 = 0.975.$$

~~$t_\alpha = 5.629$~~

$$t_\alpha = -t_{(1-\alpha)} = -t_{(0.025)}$$

~~$t_\alpha = -2.145$~~

$$Q. P(-t_{0.025} < T < t_{0.05}) = P(t_{0.975} < T < t_{0.05})$$

$$= P(T < t_{0.975}) - P(T < t_{0.05}) = P(t_{0.975}) - P(t_{0.05})$$

$$= P(T > -2.145) - P(T < 0.05) = 0.975 - 0.05$$

$$= 0.925$$

VIMP

Q. Find K s.t. $P(K < T < -1.761) = 0.045$ for a random sample of size 15 from a normal population.

→ $\mu = (15-1)\bar{x} / 30$. Using t-distribution, $\mu = 15 - 1 = \underline{\underline{14}}$.

$$P(T > K) - P(T < -1.761) = 0.045$$

$$P(T > K) - P(T < -t_{(0.05)}) = 0.045$$

$$P(T > K) - (1 - 0.05) = 0.045$$

$$\begin{aligned} P(T > K) &= 1 + 0.045 - 0.05 \\ &= 1 - 0.005 \end{aligned}$$

$$P(T > K) = 0.995$$

~~$$P(T > K) = t_K = 0.995 = t_{(1-K)}$$~~

$$\cancel{t_K} = \underline{\underline{K = -2.977}}$$

More
Easy Way

$$P(K < T < t_{0.95}) = 0.045$$

$$\alpha - \beta = 0.045$$

$$\alpha = 0.995$$

$$K = t_\alpha = -t_{1-\alpha}$$

$$K = -t_{(0.005)} = \underline{\underline{-2.977}}$$

F-Distribution:

- (i) Let, X and Y be χ^2_m and χ^2_n -distributed R.V.s, then the distribution ratio, $\frac{X/m}{Y/n} \sim F_{m,n}$ distributed, with (m, n) dof's.
- (ii) pdf is given by:-

$$f(x) = \begin{cases} \frac{\left(\frac{(n+m)}{2}\right) \left(\frac{n}{m}\right)^{\frac{n}{2}} x^{\left(\frac{n}{2}-1\right)}}{\Gamma\left(\frac{n}{2}\right) \Gamma\left(\frac{m}{2}\right) \left(1 + \frac{nx}{m}\right)^{\frac{n+m}{2}}} & , x > 0 \\ 0 & \text{otherwise} \end{cases}$$

(iii)

Q. Find $f_{0.05}$ when dof are $\begin{matrix} 6 \\ 1 \end{matrix}$ and $\begin{matrix} 10 \\ 2 \end{matrix}$.

$$\rightarrow f_{0.05}(6, 10) = \underline{\underline{3.22}}$$

R-Commands for Quantile:-

- `qchisq(p, df)`
- `qt(p, df)`
- `qf(p, df1, df2)`

Inferences :-

• Simple Random Sample:-

It is a sample in which, each voter has an equal probability of being selected in the sample, and is independently chosen from sample population.

• Parameters of Population:-

It is a numerical value that gives a characteristics of entire population. Denoted by $\underline{\bar{P}}$.

• Sample Estimates:-

Numerical Values calculated from a sample that provides estimates or approximations of population parameters.

• Statistics:-

- A fⁿ of R.V. is called statistic. It is denoted by $\underline{\underline{T(x)}}$.

- A statistic is also a R.V.

- Statistic is used to Estimate a population parameter,
i.e. $T(x)$ is an estimator of θ .
i.e. $T(x) = \hat{\theta}$

• Unbiased Estimator:-

An Estimator $T(x)$ is unbiased if $E_{\theta}(T(x)) = \theta$.

The Bias of an Estimator is

$$\text{Bias}_{\theta}(T(x)) = E_{\theta}(T(x)) - \theta.$$

An Estimator is said to be unbiased, if it's bias is zero. $\text{Bias}_\theta(T(x)) = 0$.

• Variance of $T(x)$:-

$$\underline{\text{Var}_\theta(T(x))} = E \left\{ [T(x) - E(T(x))]^2 \right\}$$

• Mean Square Error (MSE) :-

$$\underline{\text{MSE}_\theta(T(x))} = E [T(x) - \theta]^2$$

Also:

$$\underline{\text{MSE}_\theta(T(x))} = \underline{\text{Var}_\theta(T(x))} + [\text{Bias}_\theta(T(x))]^2.$$

□ Theorem :-

Let, $x = (x_1, x_2, \dots, x_n)$ be an i.i.d. sample of sample variable X with population mean μ , and Population var. σ^2 , then the A.M. i.e. $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ is an unbiased estimator of μ .

And Sample var. $s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$ is an unbiased estimator of σ^2 .

• Population Mean :- μ .

• Sample Mean :- \bar{x}

• Population Variance :- σ^2

• Sample Variance :- s^2 .

To Prove:- (i) $E(\bar{x}) = \mu$.

(ii) $E(s^2) = \sigma^2$

(i) L.H.S. = $E(\bar{x}) = \frac{1}{n} \sum$

$$= E\left(\frac{x_1 + x_2 + \dots + x_n}{n}\right) = \frac{1}{n} \cdot \left[E(x_1) + E(x_2) + \dots + E(x_n) \right]$$

$$= \frac{n\mu}{n}$$

$$= \mu.$$

L.H.S. = R.H.S.

(ii)

$$\sum_{i=1}^n (x_i - \mu)^2 = \sum_{i=1}^n [(x_i - \bar{x}) + (\bar{x} - \mu)]^2$$

$$= \sum_{i=1}^n [(x_i - \bar{x})^2 + (\bar{x} - \mu)^2 + 2(x_i - \bar{x})(\bar{x} - \mu)]$$

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

$$\sum_{i=1}^n (x_i - \mu)^2 = \sum_{i=1}^n (x_i - \bar{x})^2 + (\bar{x} - \mu)^2 n + (0).$$

Dividing Both sides by $\frac{(n-1)}$.

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

\downarrow
 s^2

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

Taking Expectation of Both sides,

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

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

Extreme
V.V.IMP.

$\left[\text{Var}(\bar{x}) = \sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} ; \sum \frac{\sigma_{x_i}^2}{n} = \sigma^2 \right]$

$$E(S^2) = \frac{n}{n-1} \times \frac{\sigma^2}{n} - \frac{n}{n-1} \cdot \sigma^2 = \sigma^2 \left[\frac{1}{n-1} - \frac{n}{n-1} \right] = 1 \times \sigma^2$$

$\therefore \underline{E(S^2) = \sigma^2}$ — (Hence, Proved).

Q. Let x_1, x_2, \dots, x_n be i.i.d. sample of size n , with population mean μ and population variance σ^2 .

Prove:- (i) $\tilde{x} = \bar{x} + 1 = \frac{\sum (x_i + 1)}{n}$ is Biased Estimator of μ .

(ii) $\tilde{s}^2 = \frac{\sum (x_i - \bar{x})^2}{n}$ is Biased Estimator of σ^2 .

\rightarrow (i) $E(\tilde{x}) = E(\bar{x} + 1) = E(\bar{x}) + E(1) = \mu + 1 \neq \mu$
 $\therefore \underline{\text{Biased Estimator}}$

$$(ii) E(\tilde{s}^2) = E \left(\frac{\sum (x_i - \bar{x})^2}{n} \times \frac{n-1}{n-1} \right) = E \left(\frac{\sum (x_i - \bar{x})^2}{n-1} \times \frac{n-1}{n} \right) = E \left(s^2 \times \frac{(n-1)}{n} \right)$$

$$= \frac{n-1}{n} E(s^2) = \left(\frac{n-1}{n} \right) \sigma^2 \neq \sigma^2 \quad \therefore \underline{\text{Biased Estimator}}$$

Hypothesis

One Sample Test:-
Data is assumed to arise as one sample from defined population.

Two Sample Test:-
Data originates in the form of two samples from two different Population.

→ Two Independent Sample Problem.

→ Two Dependent Sample Problem.

Hypothesis:- A claim that we want to test.

i) Null Hypothesis:-

Currently accepted value for a parameter.

It is opposite of Alternative Hypothesis; denoted by

H_0

ii) Alternative Hypothesis:-

Involves the claim to be tested. Denoted by H_1 / H_a .

It is formulated as deviation from the target value.

A company has stated that their straw machine makes a straws that are 4 mm diameter. A worker believes the machine no longer makes straws of the size and samples 100 straws to perform a hypothesis test with 99% confidence.

$$\rightarrow H_0: \mu = 4 \text{ mm.} \quad \text{vs straw diameter}$$

$$H_1: \mu \neq 4 \text{ mm.}$$

Doctors believe that the average teen sleeps on average less than 10 hrs. per day. A researcher believes that teens on average sleep longer. Write H_0 and H_1 .

$$\rightarrow H_0: \mu \leq 10 \text{ hrs.}$$

$$H_1: \mu > 10 \text{ hrs.}$$

One and Two Sided Tests

For an unknown population parameter θ and a fixed value θ_0 , we have:

Case	H_0	H_1	
A	$\theta = \theta_0$	$\theta \neq \theta_0$	Two sided Test
B	$\theta \geq \theta_0$	$\theta < \theta_0$	One sided test.
C	$\theta \leq \theta_0$	$\theta > \theta_0$	One sided test.

Type I and Type II Errors

Decision / Act	H_0 is True.	H_0 is not True.
H_0 is not Rejected.	Correct Decision.	Type II Error.
H_0 is Rejected.	Type I Error.	Correct Decision.

(i) Type I Error :-

The Hypothesis H_0 is True, but is Rejected.
i.e. H_1 is accepted.

(ii) Type II Error :-

The hypothesis is not rejected, although it is not True.

(a) Significance Level :-

Probability of Type I Error, i.e.

$$\underline{P(H_1 | H_0)} = \alpha$$

i.e. Probability of accepting H_1 , if H_0 is True.

(b) Level of Confidence :-

Probability of Type II Error.

$$\underline{P(H_0 | H_1)} = \beta$$

■ We try to fix α and then minimize β .

□ Power of Test :-

$$\underline{1 - \beta = P(H_1 | H_1)}$$

i.e. Prob. of making a decision in favor of the research hypothesis H_1 , if it is true. i.e.

the Prob. of detecting a correct research hypothesis.

Steps to Conduct a Statistical Test:-

Step 1:-

- Define the Assumptions for random variables of interest and specify them in terms of population Parameters. (i.e. θ or μ or σ).
- Formulate Null Hypothesis and Alternative Hypothesis.
- Fix a significance value / level i.e. α .

Step 2:- Consider a statistic, $T(x) = T(x_1, x_2, \dots, x_n)$ The distribution has to be known under Null hypothesis.

Step 3:- Consider a Critical region, K for the statistics T , i.e. a region where if T falls in this region - H_0 is rejected, such that

$$P_{H_0}(T(x) \in K) \leq \alpha$$

Step 4:-

Find $t(x) = T(x_1, x_2, \dots, x_n)$ based on sample value $x_1 = x_1, x_2 = x_2, \dots, x_n = x_n$.

Step 5:- Decision Rule:-

If $t(x)$ falls into critical region K , then null Hypothesis H_0 is rejected.

i.e. $t(x) \in K$: H_0 is rejected.
 $\Rightarrow H_1$ is significant.

If $t(x)$ falls outside K , null hypothesis is not rejected.
i.e. $t(x) \notin K$; H_0 is not rejected.
 $\Rightarrow H_0$ is accepted and is significant.

□ Tests:-

[i] One-Sample Gauss Test (z-test) :-

This test is used ~~when~~ to test whether unknown mean differs from a specific value of mean of sample.

- We assume $\sigma^2 = \sigma_0^2$ is given.

♦ Steps:-

1. (i) The random variable X follows a $N(\mu, \sigma^2)$, distribution with known variance σ^2 .

(ii) formulating H_0 and H_1 .

$H_0: \mu = \mu_0$; $H_1: \mu \neq \mu_0$ (Two sided test).

$H_0: \mu \leq \mu_0$; $H_1: \mu > \mu_0$ (One sided test).

$H_0: \mu \geq \mu_0$; $H_1: \mu < \mu_0$ (one sided test).

(iii) $\alpha = 0.05$ (used as default if not given in question).

2. Constructing a test statistics:-

We know that if X_i 's are i.i.d. then Sample Mean is Normally Distributed.

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu_0, \frac{\sigma_0^2}{n}\right)$$

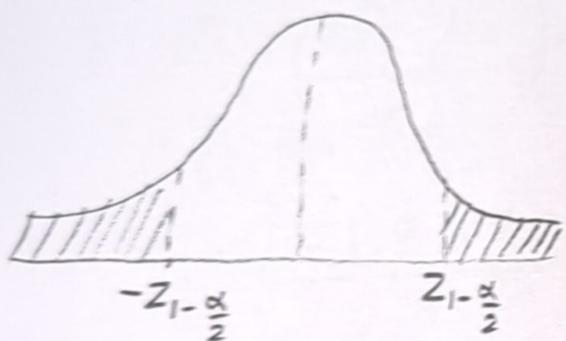
$$; T(X) = \frac{(\bar{X} - \mu_0)\sqrt{n}}{\sigma_0} \sim N(0, 1).$$

3. Critical Region.

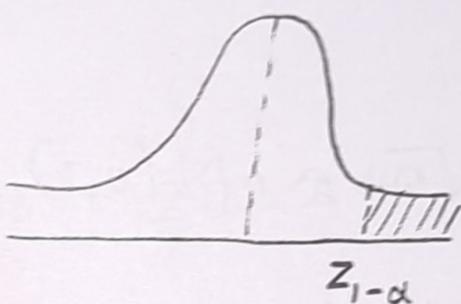
Consider the following table.

Case	H_0	H_1	Critical Region.
Two sided a.	$\mu = \mu_0$	$\mu \neq \mu_0$	$K = (-\infty, -z_{1-\frac{\alpha}{2}}) \cup (z_{1-\frac{\alpha}{2}}, \infty)$
One sided b.	$\mu \leq \mu_0$	$\mu > \mu_0$	$K = (z_{1-\alpha}, \infty)$
One sided c.	$\mu \geq \mu_0$	$\mu < \mu_0$	$K = (-\infty, z_\alpha = -z_{1-\alpha})$.

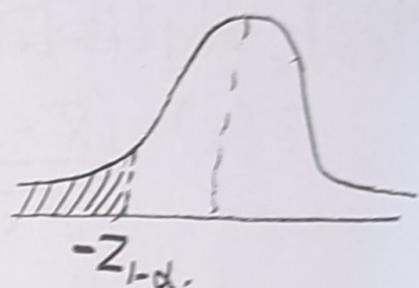
a)



b)



c)



4. Realization of Test statistics:

For observed samples x_1, x_2, \dots, x_n the arithmetic mean $\bar{x} = \frac{\sum x_i}{n}$ is used to calculate realized test statistic, $t(x)$ i.e.

$$\begin{aligned} t(x) &= T(x_1, x_2, \dots, x_n) \\ &= \frac{(\bar{x} - \mu_0)}{\sigma_0} \sqrt{n}. \end{aligned}$$

5. Decision Rule:-

If $t(x) \in K \Rightarrow H_0$ is Rejected.
 $\Rightarrow H_1$ is Accepted.

If $t(x) \notin K \Rightarrow H_0$ is Accepted.
 $\Rightarrow H_1$ is Rejected.

Q. A sample of 100 tires is taken from a lot. The mean life of tires in the sample was found to be 39350 kms with the population standard deviation of 3260 km. Test the hypothesis, at 1% level of significance that the mean life of tire is 40,000 kms.

$$\rightarrow (1) n = 100, \bar{x} = 39350 \text{ kms}, \sigma = 3260 \text{ kms}, \alpha = 0.01.$$

$$H_0: \mu = 40000 \text{ kms}$$

$$H_1: \mu \neq 40000 \text{ kms.}$$

2. Test statistic:

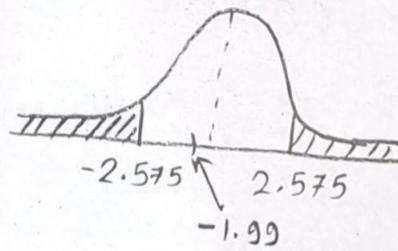
$$T(x) = \frac{(\bar{x} - \mu_0)}{\sigma_0} \sqrt{n} \sim N(0, 1).$$

3. We take third critical region.

$$K = (-\infty, -z_{1-\frac{\alpha}{2}}) \cup (z_{1-\frac{\alpha}{2}}, \infty)$$

$$= (-\infty, -z_{0.995}) \cup (z_{0.995}, \infty)$$

$$K = (-\infty, -2.575) \cup (2.575, \infty)$$



$$4. t(x) = \frac{(\bar{x} - \mu_0)}{\sigma_0} \sqrt{n} = \frac{(40000 + 39350)}{3260} \times 10 = -\underline{\underline{1.99}}$$

5. $t(x) \notin K$. $\therefore H_0$ is accepted. H_1 is rejected.

Q. A random sample of 100 recorded deaths in the US during the past year showed an average life span of 71.8 years. Assuming a population standard deviation of 8.9 years does this seem to indicate that mean life span today is greater than 70 years? Use a 0.05 level of significance.

$$\rightarrow 1. n = 100, \bar{x} = 71.8 \text{ years}, \sigma = 8.9 \text{ years}, \alpha = 0.05.$$

$$H_0: \mu \leq 70 \text{ years.}$$

$$H_1: \mu > 70 \text{ years.}$$

2. Test Statistic:-

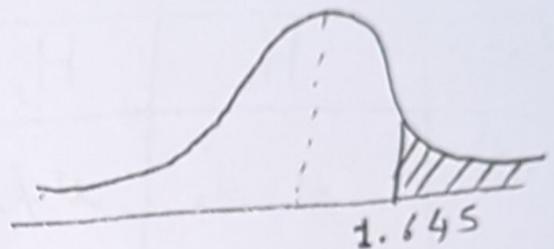
$$T(x) = \frac{(\bar{x} - \mu_0)\sqrt{n}}{\sigma_0} \sim N(0, 1).$$

3. Critical Region:-

$$K = (Z_{1-\alpha}, \infty).$$

$$K = (Z_{0.95}, \infty) = (1.645, \infty).$$

$$4. t(\alpha) = \left(\frac{\bar{x} - \mu_0}{\sigma_0} \right) \sqrt{n}$$



$$= \left(\frac{71.8 - 70}{8.9} \right) \times 10 = 2.02247.$$

5. $t(\alpha) \in K \Rightarrow H_1$ is accepted. H_0 is Rejected.

[II] t-test:-

This test is used when σ^2 is unknown.

• Sample Variance $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$.

Population Variance, $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$.

• Steps:-

1. (i)
 - (ii)
 - (iii)
- Same as Z-test.

2. Constructing a test statistic:-

$$T(x) = \frac{(\bar{x} - \mu_0)\sqrt{n}}{s_x} \text{ with } n-1 \text{ degrees of freedom.}$$

3. Critical Region:- (degree of freedom, $v = n-1$)

Case	H_0	H_1	K.
(a)	$\mu = \mu_0$	$\mu \neq \mu_0$	$(-\infty, - t_{1-\frac{\alpha}{2}}) \cup (t_{1-\frac{\alpha}{2}} , \infty)$
(b)	$\mu \geq \mu_0$	$\mu < \mu_0$	$(-\infty, - t_{1-\alpha})$
(c)	$\mu \leq \mu_0$	$\mu > \mu_0$	$(t_{1-\alpha} , \infty)$.

4. Realization of test statistics:-

For an observed sample x_1, x_2, \dots, x_n

$$\bar{x} = \frac{\sum x_i}{n} \quad \text{and} \quad t(\alpha) = \frac{(\bar{x} - \mu_0)\sqrt{n}}{s_x}$$

5. Decision Rule:

(Same as z-test).

- Q. A manufacturer of a certain brand of energy bar claims that the average saturated fat content in the bar is 0.5 gms. Will you support his claim if the 8 bars that you examined for fat content were found to contain 0.6, 0.7, 0.7, 0.3, 0.4, 0.5, 0.4 and 0.2 gms of saturated fats? Take $\alpha = 0.05$.

1. $H_0: \mu = 0.5$ gms.

$\alpha = 0.05$

$H_1: \mu \neq 0.5$ gms.

$n = 8 \quad v = 7$

$$\bar{x} = \frac{0.6 + 0.7 + 0.7 + 0.3 + 0.4 + 0.5 + 0.4 + 0.2}{8} = 0.475$$

$$2. T(x) = \left(\frac{\bar{x} - \mu_0}{S_x} \right) \sqrt{n}$$

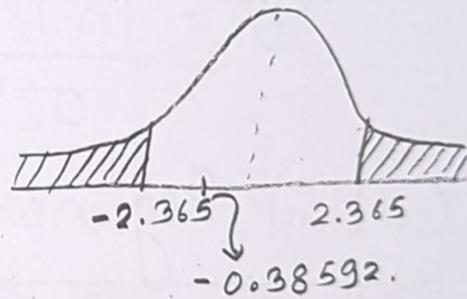
3. Critical Region:

$$K = (-\infty, -|t_{1-\frac{\alpha}{2}}|) \cup (|t_{1-\frac{\alpha}{2}}|, \infty)$$

$$= (-\infty, -|t_{0.975}|) \cup (|t_{0.975}|, \infty)$$

$$= (-\infty, -|-t_{0.025}|) \cup (|-t_{0.025}|, \infty)$$

$$K = (-\infty, -(2.365)) \cup (2.365, \infty)$$



$$4. t(\alpha) = \left(\frac{\bar{x} - \mu_0}{S_x} \right) \sqrt{n} = \left(\frac{0.475 - 0.5}{0.183225} \right) \sqrt{8} \quad S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$t(\alpha) = \underline{-0.38592}$$

5. Decision:-

$t(\alpha) \notin K \therefore H_0$ is Accepted.

H_1 is Rejected.

[III] χ^2 -Test:

- This test is used when $\sigma^2 = \sigma_0^2$ is given.
- We assume the Distribution of the Population Sample is Normal.

Steps:-

- (i) }
 (ii)
 (iii) } Same as Previous tests.

2. Constructing a test statistics:-

$$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$$

3. Critical Regions:-

Case.	H_0	H_1	K .
(a)	$\theta = \theta_0$	$\theta \neq \theta_0$	$K = (0, \chi^2_{1-\frac{\alpha}{2}}) \cup (\chi^2_{\frac{\alpha}{2}}, \infty)$
(b)	$\theta \geq \theta_0$	$\theta < \theta_0$	$K = (0, \chi^2_{1-\alpha})$
(c)	$\theta \leq \theta_0$	$\theta > \theta_0$	$K = (\chi^2_{\alpha}, \infty)$

4. Realization of test statistics:-

$$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$$

5. Decision:-

Same as Previous tests.

Q. A manufacturer of car batteries claims that the life of a company's batteries is approx. normally distributed with a standard deviation of 0.9 years. If a random sample of 10 of these batteries has a standard deviation of 1.2 years, do you think $\sigma > 0.9$ years? use $\alpha = 0.05$.

$$\rightarrow 1. H_0: \sigma \leq 0.9 \text{ years.} \quad \alpha = 0.05 \\ \sigma > 0.9 \text{ years.} \quad s = 1.2 \text{ years.} \quad n = 10. \quad v = 9.$$

$$2. \text{ Test Statistic} = \chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$$

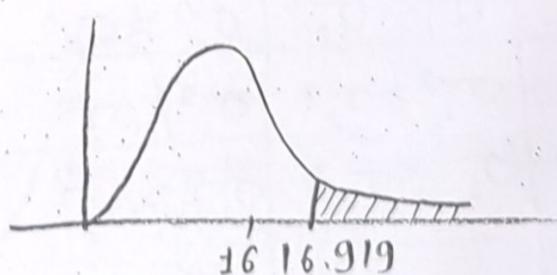
3. Critical Region:

$$K = (\chi_{\alpha}^2, \infty)$$

$$K = (\chi_{0.05}^2, \infty) = (16.919, \infty)$$

4. Realization of test statistics:-

$$\chi^2 = \frac{(10-1) \times (1.2)^2}{(0.9)^2} = 16.$$



5. Decision:-

$$\chi^2 \notin K^2$$

H_1 is Rejected. H_0 is Accepted.

[IV] F-test :-

* We use this test for comparing Variance i.e.

whether $\sigma_1^2 = \sigma_2^2$ or $\sigma_1^2 \neq \sigma_2^2$

* We assume the population to be normally Distributed.

Steps:-

1. (i) }
 (ii) }
 (iii) } → Same as Previous tests.

2. Constructing a Test Statistics:-

$$f = \frac{s_1^2}{s_2^2} \quad \text{Also, we have: } f_{\alpha}(v_1, v_2) = \frac{1}{f_{1-\alpha}(v_2, v_1)}$$

3. Critical Regions:

Case	H_0	H_1	Region (K).
(a)	$\sigma_1^2 = \sigma_2^2$	$\sigma_1^2 \neq \sigma_2^2$	$(0, f_{1-\frac{\alpha}{2}}(v_1, v_2)) \cup (f_{\alpha/2}(v_1, v_2), \infty)$
(b)	$\sigma_1^2 \geq \sigma_2^2$	$\sigma_1^2 < \sigma_2^2$	$(0, f_{1-\alpha}(v_1, v_2))$
(c)	$\sigma_1^2 \leq \sigma_2^2$	$\sigma_1^2 > \sigma_2^2$	$(f_{\alpha}(v_1, v_2), \infty)$

4. Realization of test Statistics:

$$f = \frac{s_1^2}{s_2^2}$$

5. Decision:-

Same as Previous tests.

Q. A study is conducted to compare the lengths of time required by men and women to assemble a certain product. Past experience indicates that the distribution of times for both men and women is approx. normal but variance of times for a woman is less than that of a man. A random sample of times for 11 men and 14 women gives resp. standard dev. 6.1 and 5.3. Test the hypothesis $\sigma_1^2 = \sigma_2^2$ against $\sigma_1^2 > \sigma_2^2$.

→ 1. $n_1 = 11$, $n_2 = 14$. $v_1 = 10$, $v_2 = 13$.
 $s_1 = 6.1$ $s_2 = 5.3$ $\alpha = 0.05$ - (default).

$$H_0: \sigma_1^2 \leq \sigma_2^2$$

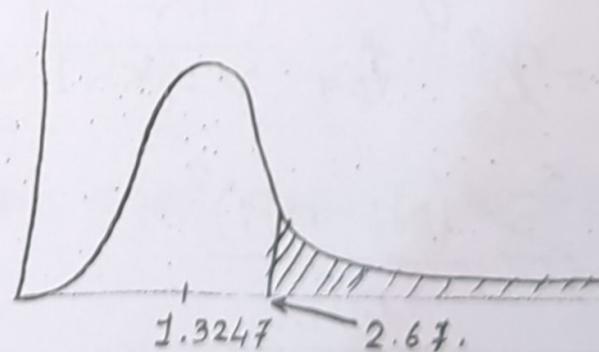
$$H_1: \sigma_1^2 > \sigma_2^2$$

$$2. f = \frac{s_1^2}{s_2^2}$$

$$3. K = (f_\alpha(v_1, v_2), \infty) = (f_{0.05}(10, 13), \infty)$$

$$K = (2.67, \infty)$$

$$4. f = \frac{(6.1)^2}{(5.3)^2} = 1.3247$$



5. $f \notin K \therefore H_1$ is Rejected.
 $\therefore H_0$ is Accepted.

□ Note:- Consider two samples X and Y in i.i.d. of size m and n from normal population, and their sample variance are given by $S_x^2 = \frac{1}{m-1} \sum (x_i - \bar{x})^2$ and $S_y^2 = \frac{1}{n-1} \sum (y_i - \bar{y})^2$
then $\left[\frac{S_x^2}{S_y^2} \sim F(m-1, n-1) \right]$

[IV] χ^2 Goodness of fit:-

- used when observed absolute frequencies are compared with the expected absolute frequencies under H_0 .

Steps :-

1. (i) }
 (ii) }
 (iii) } → Same as previous tests.

2. Test Statistic :-

$$T(x) = \chi^2 = \sum_{i=1}^k \left\{ \frac{(N_i - np_i)^2}{np_i} \right\},$$

where N_i are the absolute frequencies of observations of the sample x_i in class i , N_i is a random variable with realization n_i in the observed sample.

3. Critical Region:

$K = \chi^2_\alpha$ for $v = k-1$ degree of freedom.

$$4. \chi^2 = \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i}$$

5. If $\chi^2 > K$, then $\underline{H_0}$ is rejected.
 $\underline{H_1}$ is Accepted.

Q. Gregor Mendel conducted crossing experiment with pea plants of different shape and color. Let us look at the outcome of a pea crossing experiment with the following result:-

Crossing Result	RY	RG	EY	EG.	
Observations	315	108	101	32	$n = 556$

Mendel has hypothesis that the four different types occur in proportions of 9:3:3:1 i.e.

$$P_1 = \frac{9}{16}, P_2 = \frac{3}{16}, P_3 = \frac{3}{16}, P_4 = \frac{1}{16}.$$

$$\rightarrow 1. H_0: P(X=i) = P_i \quad ; \quad i = 1, 2, 3, 4 \quad \alpha = 0.05$$

$$H_1: P(X=i) \neq P_i \quad ; \quad v = 4 - 1 = 3 = \underline{\underline{k-1}}$$

$$2. \chi^2 = \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i} = \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i}$$

$$3. k = \chi^2_{\alpha} = 7.815$$

$$4. \chi^2 = \frac{(315 - 556 \times \frac{9}{16})^2}{556 \times \frac{9}{16}} + \frac{(108 - 556 \times \frac{3}{16})^2}{556 \times \frac{3}{16}} + \frac{(101 - 556 \times \frac{3}{16})^2}{556 \times \frac{3}{16}} + \frac{(32 - 556 \times \frac{1}{16})^2}{556 \times \frac{1}{16}}$$

$$= \frac{(2.25)^2}{34.75 \times 9} + \frac{(3.75)^2}{34.75 \times 3} + \frac{(3.25)^2}{34.75 \times 3} + \frac{(2.75)^2}{34.75}$$

$$= \frac{1}{34.75} \left[\frac{(2.25)^2}{9} + \frac{(3.75)^2}{3} + \frac{(3.25)^2}{3} + \frac{(2.75)^2}{3} \right] = \frac{16.333}{34.75}$$

$$\chi^2 = 0.47002 \quad \left. \right\} \quad 5. 0.47 < 7.815 \\ \therefore H_0 \text{ is Accepted. } H_1 \text{ is Rejected.}$$

Test Decisions Using Confidence Intervals:

1. If H_0 is Rejected at the significance level α , then there exist a $100(1-\alpha)\%$ confidence interval which yields the same conclusion as the test.
2. This is called as "Duality".

e.g. Consider 2 drugs A and B and we want to find whether average changes due to drug B is greater (higher) than for drug A i.e. $H_1 : S_B > S_A$.

So we want to decide whether H_1 is significant or not. This is equivalent to constructing $100(1-\alpha)\%$ confidence interval for $S_B - S_A$ and checking whether H_0 is significant or not.

Conditions when Tests are used :-

1. Z-test: (i) When σ is Known.
(ii) $n \geq 30$ or Data is Normally distributed.
(iii) Comparing the sample mean to a known mean of population.
2. t-Test:- (i) When σ is unknown.
(ii) $n < 30$ or Data is Approx. Normally Distributed.
3. χ^2 -Test:- When Distribution is Normal and fixed variance/ standard deviation is given.
4. F-test:- When variance need to be compared.
5. χ^2 Goodness of fit Test:- Used when observed ~~when observed~~
Absolute frequencies are compared with the expected absolute Freq. under H_0 .

Linear Regression f fitting

1. Variables (x, y, \dots)

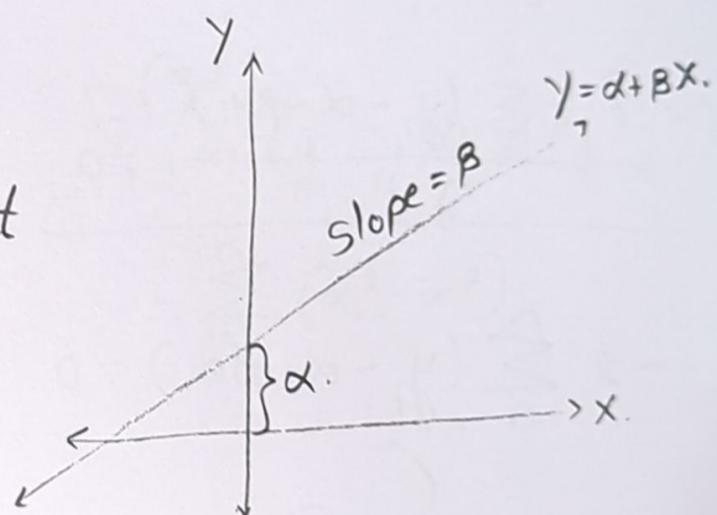
2. Regression Parameters ($\alpha, b, \alpha, \beta, \dots$)

□ Linear Model :-

$$Y = \alpha + \beta X$$

↓ ↓ ↓
 Dependent Variable (Response) y-Intercept Slope

Independent Variable (Covariate).



□ Method of Least Square :-

Consider n sets of observation, given by $P_i = (x_i, y_i)$, $i=1, 2, \dots, n$ obtained on two variables $P = (x, y)$.

The method of least squares says that, a line can be fitted to the given data set such that the errors are minimized. We need to determine α and β , denoted by $\hat{\alpha}$ and $\hat{\beta}$ s.t. the sum of squared distances bet' data points and the line $y = \alpha + \beta X$ is minimized.

$$y_1 = \alpha + \beta x_1 + e_1$$

$$y_2 = \alpha + \beta x_2 + e_2$$

$$\Rightarrow e_i = y_i - \alpha - \beta x_i$$

$$y_i = \alpha + \beta x_i + e_i$$

we have to minimize, $\sum_{i=1}^n e_i^2$

$$\therefore \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

To minimize, using principle of Maxima and Minima.

Partially Diff. w.r.t. α and β , and equating with '0'.

$$\frac{\partial (\sum_{i=1}^n e_i^2)}{\partial \alpha} = 0 ; \quad \frac{\partial (\sum_{i=1}^n e_i^2)}{\partial \beta} = 0.$$

$$\frac{\partial (\sum_{i=1}^n (y_i - \alpha - \beta x_i)^2)}{\partial \alpha} = 0 ; \quad \frac{\partial (\sum_{i=1}^n (y_i - \alpha - \beta x_i)^2)}{\partial \beta} = 0.$$

$$-2 \sum_{i=1}^n (y_i - \alpha - \beta x_i) = 0 ; \quad -2 \sum_{i=1}^n (y_i - \alpha - \beta x_i) x_i = 0.$$

Since, we are estimating $\hat{\alpha}$ and $\hat{\beta}$, replace α, β with $\hat{\alpha}, \hat{\beta}$.

$$\sum_{i=1}^n y_i - \sum_{i=1}^n \hat{\alpha} - \sum_{i=1}^n \hat{\beta} x_i = 0 . \quad | \quad \sum_{i=1}^n (x_i y_i - \hat{\alpha} x_i - \hat{\beta} x_i^2) = 0.$$

$$n\bar{y} - n\hat{\alpha} - \hat{\beta}n\bar{x} = 0.$$

$$n\hat{\alpha} = n\bar{y} - \hat{\beta}n\bar{x}$$

$$\boxed{\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}}$$

$$\sum_{i=1}^n x_i y_i - \hat{\alpha} n \bar{x} - \hat{\beta} \sum_{i=1}^n x_i^2 = 0.$$

$$\sum_{i=1}^n x_i y_i - \hat{\beta} \sum_{i=1}^n x_i^2 - n \bar{x} (\bar{y} - \hat{\beta} \bar{x}) = 0.$$

$$\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} - \hat{\beta} \left(\sum_{i=1}^n x_i^2 - n \bar{x}^2 \right) = 0.$$

$$\boxed{\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}}$$

$$\begin{aligned}
 \hat{\beta} &= \frac{\sum_{i=1}^n (x_i y_i - \bar{x}\bar{y})}{\sum_{i=1}^n (x_i^2 - \bar{x}^2)} = \frac{\sum_{i=1}^n (x_i y_i + \bar{x}\bar{y}) - \sum_{i=1}^n \bar{x}\bar{y} - \sum_{i=1}^n \bar{x}\bar{y}}{\sum_{i=1}^n (x_i^2 - \bar{x}^2)} \\
 &= \frac{\sum_{i=1}^n (x_i y_i + \bar{x}\bar{y}) - n\bar{x}\bar{y} - n\bar{x}\bar{y}}{\sum_{i=1}^n (x_i^2 - \bar{x}^2)} = \frac{\sum_{i=1}^n (x_i y_i + \bar{x}\bar{y}) - \sum_{i=1}^n \bar{x}\bar{y} - \sum_{i=1}^n \bar{x}\bar{y}}{\sum_{i=1}^n (x_i^2 - \bar{x}^2)} \\
 &= \frac{\sum_{i=1}^n (x_i y_i - \bar{x}\bar{y} + \bar{x}\bar{y} - \bar{x}\bar{y})}{\sum_{i=1}^n (x_i^2 - \bar{x}^2)} = \frac{\sum_{i=1}^n (x_i (y_i - \bar{y}) + \bar{x}(\bar{y} - \bar{y}))}{\sum_{i=1}^n (x_i^2 - \bar{x}^2)}
 \end{aligned}$$

$$\boxed{\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i^2 - \bar{x}^2)}}$$

R-Command :-

Tut-11. Q.5.

$x \leftarrow c(77, 50, \dots, 67)$

$y \leftarrow c(82, 66, \dots)$

$\text{lm}(y \sim x) \longrightarrow$ This will directly give, $\hat{\alpha}$ and $\hat{\beta}$.

Properties of Linear Regression:-

1. We should interpret $\hat{y}_i = \hat{\alpha} + \hat{\beta}x_i$; only in the interval $[x_1, x_n]$.
2. The point (\bar{x}, \bar{y}) always lies on the regression line.
3. The sum of Residuals is zero.

Residuals: It is the difference betⁿ y_i and \hat{y}_i .

Denoted by $\hat{e}_i \Rightarrow \hat{e}_i = y_i - \hat{y}_i$.

Proof:

To Prove: $\sum_{i=1}^n \hat{e}_i = 0$.

$$\begin{aligned} L.H.S. &= \sum_{i=1}^n \hat{e}_i = \sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n (y_i - (\hat{\alpha} + \hat{\beta}x_i)) \\ &= \sum_{i=1}^n (y_i - \bar{y} + \hat{\beta}\bar{x} - \hat{\beta}x_i) = \sum_{i=1}^n y_i - \sum_{i=1}^n \bar{y} + \hat{\beta}\bar{x} - \sum_{i=1}^n \hat{\beta}x_i \\ &= n\bar{y} - n\bar{y} + \hat{\beta}n\bar{x} - \hat{\beta}n\bar{x} = 0. \end{aligned}$$

L.H.S. = R.H.S.

4. The Arithmetic mean of \hat{y} is equal to the arithmetic mean of y .

To Prove: $\bar{\hat{y}} = \bar{y}$.

Proof:

$$\begin{aligned} \bar{\hat{y}} &= \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n (\hat{\alpha} + \hat{\beta}x_i) = \frac{1}{n} \sum_{i=1}^n (\bar{y} - \hat{\beta}\bar{x} + \hat{\beta}x_i) \\ &= \frac{n\bar{y} - n\hat{\beta}\bar{x} + n\hat{\beta}\bar{x}}{n} = \bar{y} \end{aligned}$$

L.H.S. = R.H.S.

$\therefore \bar{\hat{y}} = \bar{y}$

Multilinear Regression:-

When we have 2 co-variates, then the regression line becomes:

$$\hat{y} = \hat{\alpha} + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2.$$

$$\hat{y} = \hat{y}_i \quad \forall i=1,2,\dots,n.$$

$$\hat{x}_1 = \hat{x}_{1i} \quad \forall i=1,2,\dots,n. \quad \hat{x}_2 = \hat{x}_{2i} \quad \forall i=1,2,\dots,n.$$

$$y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + e_i$$

$$e_i = y_i - \alpha - \beta_1 x_{1i} - \beta_2 x_{2i}$$

We need to minimize $\sum_{i=1}^n e_i^2$.

By Using Principle of Maxima and Minima,

$$\frac{\partial \sum_{i=1}^n e_i^2}{\partial \alpha} = 0 \quad ; \quad \frac{\partial \sum_{i=1}^n e_i^2}{\partial \beta_1} = 0 \quad ; \quad \frac{\partial \sum_{i=1}^n e_i^2}{\partial \beta_2} = 0.$$

$$-2 \sum_{i=1}^n (y_i - \alpha - \beta_1 x_{1i} - \beta_2 x_{2i}) = 0$$

$$-2 \sum_{i=1}^n (y_i - \alpha - \beta_1 x_{1i} - \beta_2 x_{2i}) x_{1i} = 0.$$

$$-2 \sum_{i=1}^n (y_i - \alpha - \beta_1 x_{1i} - \beta_2 x_{2i}) x_{2i} = 0.$$

$$\Rightarrow \sum_{i=1}^n y_i - n\alpha - \beta_1 \sum_{i=1}^n x_{1i} - \beta_2 \sum_{i=1}^n x_{2i} = 0. \quad \text{---(I)}$$

$$\Rightarrow \sum_{i=1}^n y_i x_{1i} = \alpha \sum_{i=1}^n x_{1i} + \beta_1 \sum_{i=1}^n x_{1i}^2 + \beta_2 \sum_{i=1}^n x_{1i} x_{2i} \quad \text{---(II)}$$

$$\Rightarrow \sum_{i=1}^n y_i x_{2i} = \alpha \sum_{i=1}^n x_{2i} + \beta_1 \sum_{i=1}^n x_{1i} x_{2i} + \beta_2 \sum_{i=1}^n x_{2i}^2. \quad \text{---(III)}$$

Solving 3 simultaneous eq^{ns}, I, II and III,
we will get values of $\hat{\alpha}, \hat{\beta}_1, \hat{\beta}_2$.

Then Regression line will be

$$\hat{y} = \hat{\alpha} + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2.$$

When we have p covariates i.e. x_1, x_2, \dots, x_p .

then regression line will be

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p.$$

Further,

$$y_1 = \hat{\beta}_0 + \hat{\beta}_1 x_{11} + \dots + \hat{\beta}_p x_{1p} + e_1$$

$$y_2 = \hat{\beta}_0 + \hat{\beta}_1 x_{21} + \dots + \hat{\beta}_p x_{2p} + e_2$$

$$y_n = \hat{\beta}_0 + \hat{\beta}_1 x_{n1} + \dots + \hat{\beta}_p x_{np} + e_n.$$

It can be written as:

$$Y = X\beta + e$$

where,

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}_{n \times (p+1)}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}_{(p+1) \times 1} \quad e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}_{n \times 1}$$

then,

$$\hat{\beta} = (\tilde{X}' X)^{-1} X' Y$$

where, X' is Transpose of X .

R- Commands:-

T-test:

x_1, x_2, \dots, x_n are sample elements.

$X \leftarrow c(x_1, x_2, \dots, x_n).$

(i) If μ_0 is given.

`t.test(X, mu = μ_0)`

(ii) Two sided test.

`t.test(X, alternative = 'two-sided')`

(iii) One sided.