# Unit-III

## Methods of estimation:-

So far we have been discussing requisites of a good estimator now we shall briefly outline some of importants methods for obtaining such estimators.

Commonly used methods are

i) Method of maximum likihood estimation

ii) Method of minimum variance

iii) method of moments

iv) method of inverse probability

v) method of minimum chi-square.

## Maximum likelihood estimations :

The method of maximum likilihood estimation is from theoritical point of view, the most general method of estimation known is the method of maximum likihood estimators (MLE)

which was initially formulated by cf Grauss but as a general method of estimation was first introduced by prof R.A fisher and later on developed by him in a series of papers before introducing the method of we will first define likihood function.

## Likelihood Function:

Definition:

Let $x_1, x_2, \ldots x_n$ be a random sample of size n from a population with density function of the sample values $x_1, x_2, \ldots x_n$ usually denoted by $L = L(\theta)$ is their joint density function, given by

$$L = f(x_1, \theta)\, f(x_2, \theta), \cdots\cdots f(x_n, \theta) = \prod_{i=1}^{n} f(x_i, \theta) \rightarrow ①$$

L gives the relative likelihood that the random variables assume a particular set of values, $x_1, x_2, \ldots x_n$ for a given sample $x_1, x_2, \ldots x_n$. L become a function of the variable $\theta$, the parameter.

$$L = f(x_1, \theta)\, f(x_2, \theta), \cdots f(x_n, \theta) = \prod_{i=1}^{n} f(x_i, \theta) \rightarrow ①$$

It is also known as a likihood function.

The principle of maximum Likelihood consists in finding an estimator for the unknown parameter $\theta = (\theta_1, \theta_2, \cdots \theta_t)$ Say, which maximises the Likilihood Function $L(\theta)$ for variance in parameter (i.e) we wish to find $\hat{\theta} = (\hat{\theta}, \hat{\theta}, \cdots \cdot \hat{\theta}_k)$ so that

$$L(\hat{\theta}) > L(\theta) \; \forall \theta \in \textcircled{H} \; (i.e)$$

$$L(\hat{\theta}) = \sup L(\theta) \sqrt{} \; \theta \in \textcircled{H}$$

This if there exists a function $\hat{\theta} = \hat{\theta} (x_1, x_2, \cdots x_n)$ of the Sample values which maximises L for variations in $\theta$, then $\hat{\theta}$ is to be takes as an estimator of $\theta$, $\hat{\theta}$ is usually called maximum Likilihood estimation (M·L·E)

Thus $\hat{\theta}$ is the solution if any of

$$\frac{\partial L}{\partial \theta} = \theta \quad \text{and} \quad \frac{\partial^2 L}{\partial \theta^2} < 0$$

Since $L > 0$, and $\log L$ is a non-decreasing Function of L. and $\log L$ attain their extreme values (maxima or minima) of the Same value of $\hat{\theta}$; The first of the two equations can be

$$\frac{1}{L} \frac{\partial L}{\partial \theta} = \theta \Rightarrow \frac{\partial \log L}{\partial \theta} = 0$$

$$\theta = \frac{\partial \log L}{\partial \theta} = 0$$

a from which is much more convinient from particular point of view.

If $\theta$ is vector valued paramitor then $\hat{\theta} = (\hat{\theta_1}, \hat{\theta_2}, \ldots \hat{\theta_L})$ is given by the solution of simultaneous equation.

$$\frac{\partial}{\partial \theta_i} \log L = \frac{\partial}{\partial \theta_i} \log L (\theta_1, \theta_2; Dx) = 0 \quad i = 1, 2, \ldots k \rightarrow \textcircled{1}$$

The above equation ③ and ④ are usually referred to as the Likelihood equations for estimating the parameters.

Propertios of maximum Likelihood :

we make the following assumptions known as the Regularity conditions.

From theoritical point of view, the most general method of estimation Known is the method of maximum Likilihood Estimators (MLE) Which was initially formulated by CFGauss

i) The first and second order derivatives.

Viz ; $\dfrac{\partial \log L}{\partial \theta}$ and $\dfrac{\partial^2 \log L}{\partial \theta}$ exist,

and are continuous function of $\theta$ in a range R

· including . the true value $\theta_0$ of the parameter for almost

all $x$.

For every $\theta$ in R, $\left| \dfrac{\partial}{\partial \theta} \log L \right| < f(x)$ and

$\left| \dfrac{\partial^2}{\partial \theta^2} \log L \right| < f_2(x)$ where $f_1(x)$ and $f_2(x)$ are integrable

Functions over $(-\infty, \infty)$

ii) The third order derivative $\dfrac{\partial^3}{\partial \theta^3} \log L$

exists such that $\left| \dfrac{\partial^3}{\partial \theta^3} \log L \right| < M(x)$.

Where $E[m(x)] < k$ , a positive quantity

iii) For every $\theta$ in R

$E\left( \dfrac{-\partial^2}{\partial \theta^2} \log L \right) = \displaystyle\int_{-\infty}^{\infty} \left( \dfrac{-\partial^2}{\partial \theta^2} \log L \right)$

$$L \, da = J(\theta)$$

is finite and non-zero

iv) The range of integration is independent of $\theta$, But if the range of integration depends on $\theta$,

then $f(a, \theta)$ vanishes at the extrems depending on $\theta$ [This assumption is two make the difference between under the integral sign is valid.

under the above assumption MLE possess the number of important properties. which will be stated in the form of the theorem.

## Method of moments :-

This method was discovered and studied in detail by karl pearson

Let $f(x): \theta_1, \theta_2, \ldots \theta_k)$ be the density function of the parent population with $k$ parameter $\theta_1, \theta_2, \ldots \theta_k$. If $\mu_i$ denotes $r$th Moment about origin. Then

$$\mu_r' = \int_{-\infty}^{\infty} x^r f(x; \theta_1, \theta_2, \ldots \theta_k) \, dx, \quad (x = 1, 2, \ldots k) \longrightarrow ①$$

$$r = 1, 2, \ldots k$$

In general $\mu_1', \mu_2', \cdots \mu_k'$ will be function of the parameters $\theta_1, \theta_2, \cdots \theta_k$

Let $x_i$; $i = 1, 2, \cdots n$ be a random sample of size n from the given population. The method of moments consist is solving the k-equation for $\theta_1, \theta_2, \cdots \theta_n$ in terms of $\mu_1, \mu_2', \cdots \mu_k'$ and

Then replacing those moments, $\mu_r'$; $r = 1, 2, \cdots k$ by the sample moments.

example:-

$$\hat{\theta_i} = \theta_i (\hat{\mu_1'}, \hat{\mu_2'}, \cdots \mu_k') = \theta_i (\mu_1' \mu_2' \cdots \mu_k')$$

$$i = 1, 2, \cdots k$$

where $\mu_r'$ is the ith moment about origin in the sample.

Then by the method of moments $\hat{\theta_1}, \hat{\theta_2}, \cdots \hat{\theta_k}$ are the required estimators of $\theta_1, \theta_2, \cdots \theta_k$ respectively

1. a) find the maximum likilihood estimate for the parameter $\lambda$ of a poisson distribution on the basis of a sample size n also find its variance.

b) show that the sample mean $\bar{x}$, is sufficient for estimating the parameter $\lambda$ of the poisson distribution.

The probability function of the poisson distribution with parameter distribution given by

$$P(x=x) = f(x, \lambda) = \frac{e^{-\lambda} \lambda^{x}}{x!} \quad x = 0, 1, 2, \cdots$$

Likilihood function of random sample $x_1, x_2, \cdots, x_n$ of n vibrations from this population is,

$$L = \prod_{i=1}^{n} f(x_i, \lambda) = \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^{n} x_i}}{x_1! \, x_2! \cdots x_n!}$$

$$\log L = n\lambda + n\bar{x} \log L - \sum_{i=1}^{n} \log(x_i!)$$

The likelihood equation for estimating $\lambda$ is $\partial/\partial\lambda$

$$\log L = 0 => -nt \quad n\bar{x}/\lambda = 0$$

Thus the MLE for $\lambda$ is the sample mean $\bar{x}$ the variance of estimate is given by,

$$\frac{1}{v(\hat{\lambda})} = E\left\{ \frac{-\partial^2}{\partial\lambda^2} (\log L) \right\}$$

$$= E\left\{ \frac{-\partial}{\partial\lambda} (-n + n\bar{x}/\lambda) \right\}$$

$$= E\left\{-\left(-n\bar{x}/\lambda^2\right)\right\} = \frac{n}{\lambda^2} E(\bar{x})$$

$$= \frac{n}{\lambda} \left[\therefore E(\bar{x}) = \frac{n}{\lambda}\right]$$

$$v(\hat{\lambda}) = \lambda/n$$

b) For the poisson distribution with parameter $\lambda$ we have

$$\frac{\partial}{\partial \lambda} \log L = -n + n\bar{x}/\lambda$$

$$= n(\bar{x}/x^{-1})$$

$$= \phi(\bar{x}, \lambda)$$

a function of $\bar{x}$, $\lambda$ only, Hence $\bar{x}$ is sufficient for estimating $\lambda$.

3. State as precisely as possible the properties of the MLE Obtain

the MLE's of $\alpha$ and $\beta$ for a random sample from the exponential

population.

$$f(x: \alpha, \beta) = y_0 e^{-\beta(x-\alpha)}, \quad a \leq x < \infty$$
$$\beta > 0$$

and you being a constant.

The probability curve is unity

$$\therefore y_0 \int_{\alpha}^{\infty} exp\left[-\beta(x-\alpha)\right] dx$$

$$y_0 \left|e^{-\frac{\beta(x-\infty))}{-\beta}}\right|_{\alpha}^{\infty} = 1$$

$$-\frac{y_0}{\beta}(0-1) = 1 \qquad y_0 = \beta$$

$$F(x; \alpha, \beta) = \beta\, e^{-\beta(x-\alpha)} \qquad \alpha \leq x < \infty$$

if $x_1, x_2, \ldots x_n$ is a random sample of $n$ observation from this population.

$$L = \prod_{i=1}^{n} F(x_i; \alpha, \beta) = \beta \exp\left\{-\beta \sum_{i=1}^{n}(x_i - \alpha)\right\} \Rightarrow \beta^n \exp\left(-n\beta(\bar{x} - \alpha)\right)]$$

$$\log L = n \log \beta - n\beta(\bar{x} - \alpha) \longrightarrow \textcircled{1}$$

The equation $\alpha$ and $\beta$ give,

$$\frac{\partial}{\partial \alpha} \log L = 0 = n\beta \longrightarrow \textcircled{2} \text{ and}$$

$$\frac{\partial}{\partial \beta} \log L = 0 = \frac{n}{\beta} - n(\bar{x} - \alpha) \longrightarrow \textcircled{3}$$

$L$ is maximum $\Rightarrow \log L$ is maximum from $\textcircled{1}$ $\log L$ is maximum.

i.e $\hat{\alpha} = x_{(1)}$ consequently $\textcircled{3}$ gives $\frac{1}{\beta} = \bar{x} - \hat{\alpha} = \bar{x} - x_{(1)}$

$$\hat{\beta} = \frac{1}{\bar{x} - x_{(1)}}$$

Hence, MLE's for $\alpha$ and $\beta$ are given by,

$$\hat{\alpha} = x_{(1)} \text{ and } \hat{\beta} = \frac{1}{\bar{x} - x_{(1)}}$$

# Unit - IV

## Interval estimation:

interval estimation is the use of sample data to calculate an interval of possible values of an unknown population parameter; this is in contrast to point estimation, which gives a single value. Confidence interval and credible intervals

## The Forms of interval estimation:

* Confidence interval [a frequentist method]

* credible intervals [a Bayesian method]

* likelihood method intervals (a Likilihoodist method)

* Fiducial interval (a fiducial method

Other Forms of Statistical Intervals, which do not estimate parameters, include:

* tolerance intervals (an estimates of population

* prediction intervals [used mainly in regression analysis)

some other types of outcome are point estimate and decisions.

## Confidence Level

The probability that the value of parameter falls within a specified range of values

### definition:-

A confidence level refers to the percentage of all possible samples that can be expected to include the true population parameter

### For example:-

Suppose all possible samples were selected from the same population, and a confidence interval were computed for each sample.

a 95% confidence level implies that 95% of the confidence intervals would include the true population parameter. $CI = \bar{x} \pm z \dfrac{S}{\sqrt{n}}$

## Confidence Co-efficient:

The confidence co-efficient is the confidence level stated as a proportion rather than as a percentage.

eg:

if you had a confidence level of 99%, the confidence co-efficient would be .99. In general, the higher the coefficient, the more certain you are that your results are accurate.

Confidence interval:

In statistics, a confidence interval is a type of estimate computed from the statistics of the observed data. This proposes a range of palausible values for an unknown parameter. The interval has an associated confidence Level that the true parameter is in the proposed range.

$$CI = \bar{x} \pm z \frac{s}{\sqrt{n}}$$

proporti ns (or) characteristics of confidence interval:

The confidence interval (CI) is a range of values. It is expressed as a percentage and is expected to contain the best estimate of a statistical parameter.

A confidence interval of 95% mean, it is 95% certain that our population parameter lios in between this confidence interval provides a range of plausible values for a parameter.

*cluvoloped from a Sample

*for a given confidence Level

Confidence interval and confidence Limits :-

Let $x_i$ $(i=1,2,\ldots n)$ be a random sample of $n$ observations from a population involving a single unknown parameter $\theta$. Let $f(x,\theta)$ be the probability function of the parent distribution. From which the sample is drawn and Let us suppose that this distribution is continuous Let $t = t(x_1,x_2,\ldots x_n)$ a function of the sample value be an estimate of the population parameter $\theta$.

the sampling distribution given by $g(t,\theta)$ we choose once for all some value of

$\alpha$ $(5\% . \text{ or } 1\% .)$ and then

determined two constans say $c_1$ and $c_2$ such that

$$P(c_1 < \theta < c_2 | t) = 1 - \alpha$$

The quantities $c_1$ and $c_2$. so determined are known as the confidence Limits or fiducial limits and the interval $[c_1, c_2]$ within which the unknown value of the population parameter is expected to Lie, is called the confidence interval and $(1-\alpha)$ confidence co-efficient.

Thus if we take $\alpha = 0.05$ or $0.01$ we shall get $95\% .$ or $99\% .$ Confidence limits.

How to find $c_1$ and $c_2$ let $T_1$ and $T_2$ be two statistics such that,

$$P(T_1 > \theta) = \alpha_1 \; - \to ①$$

and $P(T_2 < \theta) = \alpha_2 \; - \to ②$

where $\alpha_1$ and $\alpha_2$ are constants independent of $\theta$ ① and

② can be combined to give.

$$P(T_1 < \theta < T_2) = 1 - \alpha$$

where $\alpha = \alpha_1 + \alpha_2$. Statistics $T_1$ and $T_2$ defined in ① and ② may

be taken as $c_1$ and $c_2$ defined in for example, if we take a large

sample from a normal population with mean $\mu$ and standard deviation

$\sigma$ then $z = \dfrac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$

confidence interval

$$P[-z\alpha \le z \le z\alpha] = 0.95$$

and $P[-1.96 \le z \le 1.96] = 0.95$ (from Normal probability tables)

$$\Rightarrow P\left(-1.96 \le \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \le 1.96\right) = 0.95 \Rightarrow P\left[\bar{x} - 1.96\frac{\sigma}{\sqrt{n}} \le \mu \le \bar{x} + 1.96\frac{\sigma}{\sqrt{n}}\right]$$

Thus,

$\therefore \left(\bar{x} \pm 1.96 \dfrac{\sigma}{\sqrt{n}}\right)$ are 95% confidence limits for the unknown

parameter $\mu$, the population mean and the interval

$\left[\bar{x} - 1.96\,\sigma/\sqrt{n}, \; \bar{x} + 1.96\,\sigma/\sqrt{n}\right)$ is called the 95% confidence interval

Also

$$P(-2.58 \leq z \leq 2.58) = 0.99 \text{ or}$$

$$P\left(-2.58 \leq \frac{\bar{x}-\mu}{\sigma/\sqrt{n}} \leq 2.58\right) = 0.99$$

$$P\left(\bar{x}-2.58\frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x}+2.58\,\sigma/\sqrt{n}\right) = 0.99$$

Hence 99% confidence limits for $\mu$ are;

$$\bar{x} \pm 2.58\,\sigma/\sqrt{n} \quad \text{and} \quad 99\%$$

confidence interval for $\mu$ is

$$\left(\bar{x}-2.58\,\sigma/\sqrt{n}\,,\ \bar{x}+2.58\,\sigma/\sqrt{n}\right)$$

## Test of significance for single mean

we have proved that if $x_i = (i=1,2,\cdots n)$ is a random sample of size n from a normal population with mean $\mu$ and variance $\sigma^2$ then the Sample mean is distributed normally with mean $\mu$ and variance $\sigma^2/n$ that is $\bar{x}$ Follows normal $\mu, \sigma^2/n$

that is $\bar{x} \sim N(\mu, \sigma^2/n)$ even in random sampling from non normal population provided the sample size n is Large (c.f)

thus for large samples the S.D normal variate is corresponding to $\bar{x}$ is $z = \dfrac{\bar{x}-\mu}{\sigma/\sqrt{n}}$

under the null hypothesis $H_0$. That the sample has been drawn From the population with mean $\mu$ and variance $\sigma^2$ that is no significant difference and the sample mean $(\bar{x})$ and population mean $(\mu)$ The test statistics for large sample is $z = \dfrac{\bar{x}-\mu}{\sigma}$, $\sim N(0,1)$

## Remarks:

i) if the population S.D $\sigma$ is unknown then we use in estimate provide by the sample variance given by,

$$\hat{\sigma^2} = s^2, \quad \hat{\sigma} = s \quad (\text{for large Samples})$$

ii) confidence limits for $\mu$ 1.95% confidence interval for $\mu$ is given by.

$$|z| \leq 1.96, \text{ that } \left|\dfrac{\bar{x}-\mu}{\sigma/\sqrt{n}}\right| \leq 1.96 \text{ which}$$

(i.e)

$$\bar{x} - 1.96\left(\sigma/\sqrt{n}\right) \leq \mu \leq \bar{x} + 1.96\left(\sigma/\sqrt{n}\right)$$

iii) The confidence limits for any parameter $(p, \mu \ldots)$ are also known as its fiducial limits problem

1. A sample of 900 member has a mean 3.4 cm and S.D 2.61 cm is the sample from a large population of mean 3.25 cm and S.D 2.61 cm if the population is normal and its mean is unknown, find the 95% and 98% fiducial limits of true mean.

Null hypothesis:

$$H_0 : \mu = 3.25$$

Alternative hypothesis (H1): $\mu \neq 3.25$

(two tailed) test statistics under Ho. the test statistics is

Test Statistics:

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim N(0,1) \text{ since } n \text{ is large}$$

$\bar{x} = 3.4$ cm $n = 900$ $\mu = 3.25$ cm $\sigma = 2.61$

$$Z = \frac{3.4 - 3.25}{\frac{2.61}{\sqrt{900}}}$$

$$= 0.15 \times 30/2.61$$

$$|Z| = 1.73$$

Conclusion:

since $|Z| < 1.96$ we conclude that the data don't provide is any evidence against the null hypothesis (Ho) which may.

Therefore be accepted at 5% Level of significance 95%.

Fiducial limits for the population

Test of significance difference of mean :

Let $\bar{x}$ be the mean of a sample of size $n$ , from a population with mean $\mu$ , and variance $\sigma^2$, and Let $\bar{x_2}$ be the mean of and indipendent random sample of sizes $n_2$ from another population with $\mu_2$ and variance $\sigma_2^2$. then since samplo sizes are large,

$$\bar{x}_1 \sim N\left(\mu_1, \frac{\sigma^2}{n_1}\right) \text{ and } \bar{x_2} \sim N\left(\mu_2, \frac{\sigma_2^2}{n_2}\right)$$

also $\bar{x}_1 - \bar{x_2}$ , being the difference of two independent normal variates is also a normal variate. The value of $z$ (S.N.V) corressponding $\bar{x}_1 - \bar{x_2}$ is given by,

$$z = \frac{(\bar{x}_1 - \bar{x_2}) - E(\bar{x}_1 - \bar{x_2})}{V(\bar{x}_1 - \bar{x_2})} \sim N(0,1)$$

under the null hypothesis . $H_0: \mu_1 = \mu_2$ , that there is no significant difference between the sample means, we get

$$E(\bar{x}_1 - \bar{x_2}) = E(\bar{x}_1) - E(\bar{x_2})$$

$$= \mu_1 - \mu_2 = 0$$

$$V(\bar{x}_1 - \bar{x_2}) = V(\bar{x}_1) - V(\bar{x_2})$$

The co-variance term vanishes, since the Sample means $\bar{x_1}$ and $\bar{x_2}$ are independent there under $H_0 : \mu_1 = \mu_2$. The test statistics becomes (for large Samples)

$$z = \frac{(\bar{x_1} - \bar{x_2})}{\sqrt{\sigma_1^2/n_1 + \frac{\sigma_2^2}{n_2}}} \sim N(0,1)$$

$$z = \frac{\bar{x_1} - \bar{x_2}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0,1)$$

Remarks :

1. If $\sigma_1^2 = \sigma_2^2 = \sigma^2$ that if the samples have been drawn from the population with common S.D $\sigma$ then under $H_0 : \mu_1 = \mu_2$

$$z = \frac{\bar{x_1} - \bar{x_2}}{\sigma\sqrt{1/n_1 + 1/n_2}} \sim N(0,1)$$

2. If in $z = \dfrac{\bar{x_1} - \bar{x_2}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0,1)$, $\sigma$ is not known it is based

on Sample variance.

## Unit-$\overline{V}$

### Students 't' distribution :

Let $a_i$ $(i=1,2,\dots n)$ be a random sample of size $n$ from a normal population with mean $\mu$ and variance $\sigma^2$ then students $t$ is defined by the statistics.

$$\boxed{t = \frac{\bar{x}-\mu}{\delta/\sqrt{n}} \quad \longrightarrow \text{①}}$$

where $\bar{x} = \frac{1}{n}\sum\limits_{i=1}^{n} a_i$ is the sample mean and $s^2 = \frac{1}{n-1}$

$\sum\limits_{i=1}^{n} (x_i-\bar{x})^2$ is an unbiased estimate of the population variance $\sigma^2$

and it follows student $t$ distribution with $V=(n-1)$ d.f with probability density function.

$$f(t) = \frac{1}{\sqrt{v}B\left(\frac{1}{2},\frac{v}{2}\right)} \cdot \frac{1}{\left(1+\frac{t^2}{v}\right)\left(\frac{v+1}{2}\right)} \quad -\infty < t < \infty \longrightarrow \text{②}$$

### Derivation of students t distribution :

The expression ① can be re written as

$$t^2 = \frac{n(\bar{x}-\mu)^2}{s^2}$$

$$= \frac{n(\bar{x}-\mu)^2}{ns^2/(n-1)}$$

$$\Rightarrow \frac{t^2}{(n-1)} = \frac{(\bar{x}-\mu)^2}{\sigma^2/n} \cdot \frac{1}{ns^2/\sigma^2}$$

$$= \frac{(\bar{x}-\mu)^2 (\sigma^2/n)}{ns^2/\sigma^2}$$

Since $x_i (i = 1, 2, \ldots n)$ is a random sample, mean $\mu$ variance $\sigma^2$

$$\bar{X} \sim N(\mu, \sigma^2/n)$$

$$\Rightarrow \frac{\bar{x}-\mu}{\sigma/\sqrt{n}} \sim N(0,1)$$

Hence, $\dfrac{(\bar{x}-\mu)^2}{\sigma^2/n}$ being the square

of a standard normal variate is a chi-square variate with 1 d.f

$\dfrac{ns^2}{\sigma^2}$ is a $\chi^2$ variate with $(n-1)$ d.f

$$dF(t) = \frac{1}{B(\frac{1}{2}, \frac{v}{2})} \cdot \frac{(t^2/v)^{\frac{1}{2}-1}}{\left(1 + \frac{t^2}{v}\right)^{(v+1)/2}} \, d\left(\frac{t^2}{v}\right) \qquad 0 < t^2 < \infty$$

$$\therefore \text{where } v = (n-1)]$$

$$= \frac{1}{\sqrt{v} \, B\left(\frac{1}{2} \cdot \frac{1}{2}\right)} \cdot \frac{1}{\left(1 + \frac{t^2}{v}\right)\left(\frac{v+1}{2}\right)} \, dt \qquad -\infty < t < \infty$$

the factor 2 disappearing since the integral from $-\infty$ to $\infty$ must be unity.

## Application of t distribution :-

Than t distribution has a wide number of application in statistics Some of which are annumeratad below,

i) To test if the Sample mean ($\bar{x}$) differs, significantly from the hypothetical value $\mu$ of the population mean.

ii) To test the significance of the difference between two Sample means.

iii) To test the significance of an observed Sample correlation co-efficient and Sample regression co-efficient.

In the following Sections we will discuss these application in detail One by one.

## t - Test For single mean :

i) If a random Sample $x : (i=1,2,\cdots n)$ of size n has been drawn from a normal population with a specified mean Say $\mu_0$ or

ii) if the sample mean differs significantly from the hypothetical Value $\mu_0$.

under the Null hypothesis : $H_0$

test statistics :

$$t = \frac{\bar{x} - \mu_0}{S/\sqrt{n}}$$

where $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$ and

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

Follows the students t - distribution with $(n-1)$ d.f we now compare the calculate value of t with the tabulated value of certain Level of significance.

If calculated $|t| > $ tabulated t. null hypothesis is rejected and if calculated $|t| < $ tabulated t. Ho may be accepted at the level of Significance adopted.

## t- Test for difference of mean :

Suppose we want to test IF two independent samples $x_i(i=1,2,\dots n)$ and $y_j(j=1,2,\dots n)$ of sizes $n_1$ and $n_2$ have been drawn from two normal populations with mean $\mu x$ and $\mu y$ respectively.

under the null hypothesis Ho, that the samples have been drawn from the normal population with mean $\mu x$ and $\mu y$ and under the assumption that the population variance are equal i.e

$$\sigma x^2 = \sigma y^2 = \sigma^2 \text{ say}$$

$$t = \frac{(\bar{x} - \bar{y}) - (\mu x - \mu y)}{s\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

where $\bar{x} = \dfrac{1}{n_1} \displaystyle\sum_{i=1}^{n} x_i$

$\bar{y} = \dfrac{1}{n_2} \displaystyle\sum_{j=1}^{n} y_i$

$S^2 = \dfrac{1}{n_1+n_2-2} \displaystyle\sum_i \left[ (x_i-\bar{x})^2 + \sum_j (y_j-\bar{y})^2 \right]$

is an unbiased estimate of the common population variance $\sigma^2$
Follow student's t-distribution with $(n_1+n_2-2)$ d.f

proof:

$Z = \dfrac{(\bar{x}-\bar{y}) - E(\bar{x}-\bar{y})}{\sqrt{V(\bar{x}-\bar{y})}} \sim N(0,1)$

$E(\bar{x}-\bar{y}) = E(\bar{x}) - E(\bar{y})$

$V(\bar{x}-\bar{y}) = V(\bar{x}) + V(\bar{y})$

$\qquad = \dfrac{\sigma x^2}{n_1} + \dfrac{\sigma y^2}{n_2} \Rightarrow \sigma^2 \left( \dfrac{1}{n_1} + \dfrac{1}{n_2} \right)$

$\qquad = \dfrac{(\bar{x}-\bar{y}) - (\mu_x - \mu_y)}{\sqrt{\sigma^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim N(0,1) \longrightarrow \text{①}$

let

$\chi^2 = \dfrac{1}{\sigma^2} \left[ \displaystyle\sum_{i=1}^{n} (x_i-\bar{x})^2 + \displaystyle\sum_{j=1}^{n_2} (y_j-\bar{y})^2 \right]$

$\qquad = \left[ \displaystyle\sum_i (x_i-\bar{x})^2/\sigma^2 \right] + \left[ \displaystyle\sum_j (y_j-\bar{y})^2/\sigma^2 \right]$

$\qquad = \dfrac{n_1 s x^2}{\sigma^2} + \dfrac{n_2 s y^2}{\sigma^2} \longrightarrow \text{②}$

Since,

$$\frac{n_1 s_x^2}{\sigma^2} \text{ and } \frac{n_2 s_x^2}{\sigma^2}$$

Variables with $(n_1-1)$ and $(n_2-1)$ d.f respectively by the additive property of chi-square defined in ② is a $\chi^2$ variate with $(n_1-1) + (n_2-1)$ i.e $n_1+n_2-2$ d.f

Further since sample mean and sample variance are independently distributed $\xi$ and $\chi^2$ are independent random variable.

$$t = \frac{\xi}{\sqrt{\dfrac{\chi^2}{n_1+n_2-2}}}$$

$$= \frac{(\bar{x}-\bar{y})-(\mu_x-\mu_y)}{\sigma^2 \sqrt{1/n_1 + 1/n_2}} \times \frac{1}{\dfrac{1}{n_1+n_2-2}\left\{ \sum_i (x_i-\bar{x})^2 + \sum_j (y_i-\bar{y})^2 \right\}/\sigma^2}$$

$$= \frac{(\bar{x}-\bar{y})-(\mu_x-\mu_y)}{S \sqrt{1/n_1 + 1/n_2}}$$

$$S^2 = \frac{1}{n_1+n_2-2}\left[ \sum_i (x_i-\bar{x})^2 + \sum_j (y_i-\bar{y})^2 \right]$$

and it follows students t distribution with $(n_1+n_2-2)$ d.f

# F-distribution:

## Definition:

If x and y are two indipendent chi-square variates with $v_1$ and $v_2$ d.f respectively them.

F statistics is defined by

$$F = \frac{X/v_1}{Y/v_2}$$

## F test for Equality of two population variances:

Suppose we want to test ii) whether two indipendent samples. $x_i$ (i=1,2,...n), $y_j$ (j=1,2,...n) have been drawn from the normal populations with same variance $\sigma^2$ (say) or

i) whether the two independent estimates of the population variance are homogenious are not.

null hypothesis: Ho

i) $\sigma x^2 = \sigma y^2 = \sigma^2$

$$F = \frac{Sx^2}{Sy^2}$$

where.

$$Sx^2 = \frac{1}{n_1 - 1} \sum_{i=j}^{n} (x_i - \bar{x})^2 \text{ and}$$

$$Sy^2 = \frac{1}{n^2-1} \sum_{J=1}^{n_2} (y_j - \bar{y})^2$$

proof:

$$F = \frac{Sx^2}{Sy^2} \left[ \frac{n_1}{n_1-1} Sx^2 \right] / \left[ \frac{n_2}{n_2-1} Sy^2 \right]$$

$$= \left[ \frac{n_1 Sx^2}{\sigma x^2} \cdot \frac{1}{(n_1-1)} \right] / \left[ \frac{n_2 Sy^2}{\sigma y^2} \cdot \frac{1}{n_2-1} \right]$$

$$\left[ \because \sigma x^2 = \sigma y^2 = \sigma^2 \right] \text{under } H_0 \cdots n_2]$$

Since $\dfrac{n_1 Sx^2}{\sigma x^2}$ and $\dfrac{n_2 Sy^2}{\sigma y^2}$ are independent chi-square

with $(n_1-1)$ and $(n_2-1)$ df respectively, F follows Snedecor's

F distribution with $(n_1-1), (n_2-1)$ d.f.

Constants of F distribution:

$$\mu_r' \text{(about origin)} = E(F^r)$$

$$= \int_0^\infty F^r f(F) \, dF$$

$$= \frac{(\nu_1/\nu_2)^{\nu_1/2}}{B\left(\nu_1/2, \nu_2/2\right)} \int_0^\infty F^r \frac{f^{(\nu_1/2)-1}}{\left(1+\nu_1/\nu_2 F\right)^{\nu_1+\nu_2/2}} dF \longrightarrow ①$$

To evaluate the integrals put $= v_1/v_2 \, f = y$ so that $df = v_2/v_2 \, dy$

$$\mu_r' = \frac{[v_1/v_2]^{u/2}}{B(u_1/2, v_2/2)} \int_0^\infty \frac{(v_2/v_1, 9)^{r+(v_1/2)-1}}{(1+y)(v_1+v_2)/2} \cdot \left[\frac{v_2}{v_1}\right] y$$

$$= \frac{\left(\frac{v_2}{v_1}\right)^r}{B(u/2, v_2/2)} \int_0^\infty \frac{y^{r+(u/2)-1}}{(1+9)[(u/2)+r+(v_2/2)-r]} \, dy$$

$$\Rightarrow \left[\frac{v_2}{v_1}\right]^r \cdot \frac{1}{B(u/2, v_2/2)} \, B\left[r+\frac{v_1}{2}, \frac{v_2}{2}-r\right]$$