

Unit-IV

Sampling and Statistical Inference.

Sampling is a technique of selecting a subset of a population to make statistical inference from it and estimate characteristics of the entire population.

Eg:- exit poll, vaccine trials.

→ Defns:

1) Population:

It is the set of all aggregate objects under study. (denoted by N)

2) Sample:

It is the subset of the population that is considered to study the behaviour of the population. (sample size denoted by n)

If $n \geq 30$, it is considered to be large sample.

- * The statistical constants of popn such as mean, variance etc are called parameters.
- * The statistical constants of samples which are used to estimate the parameters are called statistics. Eg: \bar{x} , s^2

Sampling can be done in diff. ways but it is important to have a random

Sample

- With replacement
- Without replacement

→ Sampling Distribut^n

It is the frequency distribut^n of a statistic over many random samples from a single population.

If we have diff. samples of size n from a population of size N , for each of these samples we can compute mean, variance etc. They will not be same.

If we grp these acc. to their freqs. then the freq distribut^n so generated is called the sampling distribut^n. In particular, we can have sampling distribut^n of mean, var. etc. The std. deviation of the sampling distribut^n is called the std. error.

→ Questions

- A popl^n consists of nos $\{1, 2, 3\}$ form a sampling distribut^n of the mean for random sample of size 2 with replacement.

$$N=3, \mu = a, \sigma^2 = \frac{1}{N} \sum (x_i - \mu)^2$$

$$\sigma^2 = \frac{1}{3} \{1+1\}$$

$$\sigma^2 = 2/3$$

$$S = \{(1,1), (1,2), (1,3), (2,1), (2,2), (2,3), (3,1), (3,2), (3,3)\}$$

$$\text{Means} = \{1, 1.5, 2, 1.5, 2, 2.5, 2, 2.5, 3\}$$

x_i	$f(x_i)$	$n_i f_i$	$x_i \bar{x}$	$(x_i - \bar{x})^2$	$f_i (x_i - \bar{x})^2$
1	1	1	-1	1	1
1.5	2	3	-0.5	0.25	0.5
2	3	6	0	0	0
2.5	2	5	0.5	0.25	0.5
3	1	3	1	1	1
		18			3

$$\text{Mean}, \bar{x} = \frac{\sum x_i f_i}{\sum f_i} = \frac{18}{9} = 2.$$

$$\text{Variance } \sigma^2 = \frac{1}{\sum f_i} \sum f_i (x_i - \bar{x})^2$$

$$\sigma^2 = \frac{1}{9} \times 3 = \frac{1}{3}$$

$$\sigma^2 = 0.33$$

$$\sigma^2 = \frac{\sigma^2}{n} = \frac{2}{3} \times \frac{1}{2}$$

⇒ Central Limit Theorem

In a popl^n N with mean μ & variance σ^2 , if we take sufficiently large random samples with replacement, then the dist^n of the sample means will be normally distributed.

→ Defⁿ: (a) (c) (d) (e) (f)

i) Hypothesis: The assumⁿ that we make regarding the parameters of the populatⁿ.

ii) Test of significance/hypothesis

An imp aspect of sampling theory is to assert that the parameter of the poplⁿ is the same as the statistic obtained of the random process. The process that decides whether to accept the hypothesis or not is called test of significance.

In order to arrive at a decision, we have to make certain assumptions known as hypotheses. The hypo is accepted or rejected based on some statistical test. We decide whether the sample statistic is significantly different from the parameter at a desired level of significance. Hence, these tests are called, tests of sign

There are 2 types of hypotheses.

i) Null hypothesis

The hypo which assumes that there is no significant difference b/w sample statistic and the parameter of the populatⁿ.

It is denoted by H_0

ii) Alternate hypo

Under this hypo, we assume that there

is a significant difference b/w sample & populatⁿ statistics.

→ Type 1 and Type 2 errors

Type 1: When a true hypo is rejected, it is called type 1 error.

Type 2: When a false hypo is accepted, it is called type 2 error.

Hypothesis	Decision	
	Accepted	Rejected
True	✓	Type 1
False	Type 2	✓

→ Level of Significance (los)

The prob level below which we reject a hypo is called the level of significance.

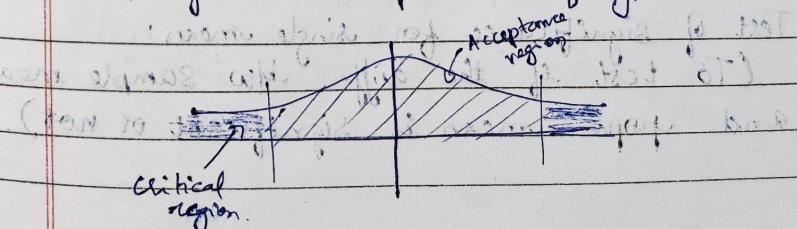
It can be 1%, 2%, 5%, 10% etc.

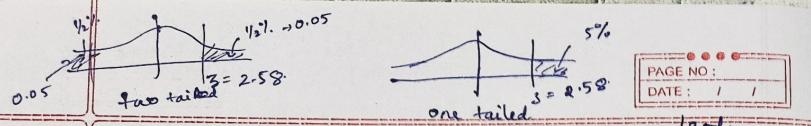
In general, 1% or 5% los is considered

→ Critical region & Acceptance regions.

C-region is the regions which corresponds to the rejectⁿ of hypo.

A-region → acceptance of hypo.





→ One tailed and two tailed ~~hypothesis~~ tests

A test on statistical hypotheses where the alternate hyno is one sided is called one tailed test.

→ Confidence Interval.

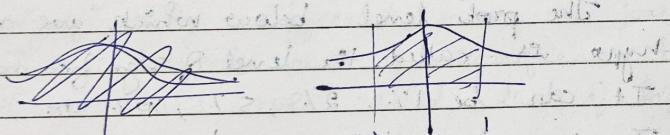
Let μ_s and σ_s be the mean & S.D of the sampling distribⁿ of a test statistic S. If the sampling dist. of S is normal then S lies in $(\mu_s - \sigma_s, \mu_s + \sigma_s)$ what % of times.

$$z = \frac{x - \mu}{\sigma}$$

$$z_1 = \frac{\mu - \mu - \mu}{\sigma}$$

$$z_2 = \frac{\mu + \sigma - \mu}{\sigma}$$

$$z_1 = -1.96 \text{ and } z_2 = 1.96 \text{ at } 5\%$$



$$A(3) = 0.3413 + 0.3413$$

$$= 68.26\%$$

→ Test For large Samples. (z-test)

- Test of significance for single mean:

(To test if the diff b/w sample mean and poplⁿ mean is significant or not)

Suppose we have a poplⁿ of size N with mean μ and variance σ^2 . Let us draw a sample of size n from it. Let \bar{x} be the sample mean. The std. error of mean of a random sample of size n is given by;

$$\text{S.E. of mean} = \frac{\sigma}{\sqrt{n}}$$

- Under the null hypothesis, there is no diff b/w std mean & sample mean.

$$H_0: \bar{x} = \mu$$

$$\text{The test statistic is } z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

- If σ is not known, the test statistic,

$$z = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

where s is the S.D. of the sample.

- If l.o.s.e is δ and z_α is the critical value, then

$$-z_\alpha \leq |z| = \left| \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \right| \leq z_\alpha$$

- The limits of poplⁿ mean, μ are given by

$$\bar{x} - z_\alpha \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + z_\alpha \frac{\sigma}{\sqrt{n}}$$

At 5% l.o.s., 95% confidence limits are

$$\bar{x} - 1.96 \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + 1.96 \frac{\sigma}{\sqrt{n}}$$

Always remember to take care of signs

At 1% l.o.s., 99% confidence limits are

$$\bar{x} - 2.58 \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + 2.58 \frac{\sigma}{\sqrt{n}}$$

\rightarrow Problems :-

1. A normal poplⁿ has a mean of 6.8 and S.D. of 1.5. If sample of 400 members gave a mean of 6.75. Is the difference significant.

Null hypothesis, $H_0 : \bar{x} = \mu \rightarrow \text{JC}$

$$\mu = 6.8, \sigma = 1.5, n = 400, \bar{x} = 6.75$$

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{6.75 - 6.8}{1.5/20} = -0.2$$

$$= -0.2$$

$|z| < z_{\alpha}$ for 5% & 1% los

$\Rightarrow H_0$ is accepted

(*) if $\bar{x} = 6.95, z = 2 > 1.96$

$|z| > z_{\alpha}$ for 5% los

H_0 is rejected

and as $|z| < z_{\alpha}$ for 1% l.o.s. ($z < 2.58$)

H_0 is accepted.

2. The mean weight for a random sample of size 100 is 64 g. The S.D. of the weight dist. of poplⁿ is 3g. Can we say that the mean weight of the poplⁿ is 66 g is at 5% los? Also set up 95% and 99% confidence interval for the mean weight of the population.

$$n = 100, \bar{x} = 64, \sigma = 3, \mu = 66$$

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{64 - 66}{3/10} = -6.6667$$

$$-6.6667 < -1.96$$

$|z| > z_{\alpha}$ for 5% & 1% los.

$\Rightarrow H_0$ is not accepted.

for 95% confidence interval, $|z| < 1.96$

$$1.96 > \frac{64 - \mu}{3/10}$$

$$0.588 > 64 - \mu$$

$$\mu > 64 - 0.588$$

$$\mu > 63.412$$

95% confidence limits are

$$64 - 1.96 \times 0.3 < \mu < 64 + 1.96 \times 0.3$$

$$\text{approx. with } 63.412 < \mu < 64.588$$

99% confidence limits are

$$64 - 2.58 \times 0.3 < \mu < 64 + 2.58 \times 0.3$$

$$63.22 < \mu < 64.774$$

3. $n = 500$, $\bar{x} = 3.4$, $\mu = 3.25$
 $s = 1.71$

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{3.4 - 3.25}{1.71/\sqrt{500}} = 8$$

$$= 0.15 \times 2.3607$$

1.71

$$1.96 < = 1.9614 < 2.58$$

5% error margin is off 1%

H_0 : accepted at 1% los.
 rejected at 5% los

$$| \bar{x}_1 - \bar{x}_2 | < 2.58$$

→ Test for Difference b/w two means

Let \bar{x}_1, \bar{x}_2 be the means of 2 samples
 of sizes n_1 and n_2 , from population
 with means μ_1 and μ_2 and S.D. is
 σ_1 and σ_2 . As the two samples are

independent for large values of n_1 and n_2 , $\bar{x}_1 - \bar{x}_2$ is approximately normally distributed and then,

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \text{ is a std. normal variable.}$$

Under the Null Hypothesis:

$$\mu_1 = \mu_2, \text{ then } z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

① Note:

- i) If σ_1^2 and σ_2^2 are unknown, then we can use the sample SD's ratio

$$z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

- ii) If the 2 populations have a common variance σ^2 , then

$$z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2}}}$$

→ Questions.

- 1) Two random samples of 100 students from two schools A and B were drawn. The CGPA of students from A had mean 2.82 with SD. 0.63 and from

school B had mean 2.43 with s.d. 0.65. Does the data indicate any difference in the mean CGPA of students from the schools?

$$\Rightarrow n_1 = n_2 = 100, \bar{x}_1 = 2.82, s_1 = 0.63 \\ \bar{x}_2 = 2.43, s_2 = 0.65$$

$$z = \frac{2.82 - 2.43}{\sqrt{\frac{0.63^2}{100} + \frac{0.65^2}{100}}} = \frac{0.39}{\sqrt{0.3969 + 0.4225}} = 0.39$$

$$z = \frac{39}{0.9052} = 43.084$$

$z > z_{\alpha}$ for 1% level. $\therefore H_0$ is rejected.

- 2) The mean heights in 2 large samples of size 1000 and 2000 are 67.5 inches and 68 inches resp. Can the samples be regarded as drawn from normal popl' with s.d. 2.5 inches.

$$n_1 = 1000, \bar{x}_1 = 67.5, \sigma_1 = 2.5$$

$$n_2 = 2000, \bar{x}_2 = 68.0$$

$$H_0: \mu_1 = \mu_2 \text{ i.e. no difference}$$

$$z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} = \frac{67.5 - 68}{\sqrt{\frac{1}{1000} + \frac{1}{2000}}} = \frac{-0.5}{\sqrt{\frac{1}{1000} + \frac{1}{2000}}} = -0.5$$

$$z = \frac{0.5}{\sqrt{\frac{1000}{2 \times 10^6}}} = \frac{0.5}{\sqrt{0.0005}} = \pm 1$$

$|z| > z_{\alpha}$ for 1% & 5%.

$\therefore H_0$ is Rejected.

$$3) n_1 = n_2 = 100, \bar{x}_1 = 1400 \text{ hrs}, s_1 = 200 \text{ hrs} \\ \text{and } n_2 = 200 \text{ hrs. } \bar{x}_2 = 1200 \text{ hrs. } s_2 = 100 \text{ hrs.}$$

$$z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{1400 - 1200}{\sqrt{\frac{200^2}{100} + \frac{100^2}{100}}} = \frac{200}{\sqrt{500}} = 2.23607$$

$$= \frac{200}{\sqrt{400+100}} = \frac{200}{\sqrt{500}} = 2.23607$$

$= 2.23607 > z_{\alpha}$ for 1% & 5%.

$\therefore H_0$ is Rejected.

→ Test of Significance for Proportion.

- For Single Proportion.

$$Z = \frac{(x - \mu)}{\sigma} = \frac{x - np}{\sqrt{npq}}$$

where,

$x \rightarrow$ no. of successes in independent trials

$p \rightarrow$ prob. of success in each trial

$$q = 1 - p$$

$$\sigma^2 = pq$$

Questions

- 1) A die is thrown 9000 times and a three or 3 or 4 is observed 3240 times. Can the die be regarded as unbiased.

$$H_0: \text{die is unbiased.}$$

$$n = 9000, p = \frac{2}{6} = \frac{1}{3}, q = \frac{2}{3}$$

$$x = 3240$$

$$Z = \frac{3240 - 9000 \cdot \frac{1}{3}}{\sqrt{9000 \cdot \frac{1}{3} \cdot \frac{2}{3}}}$$

$$= \frac{240}{\sqrt{2000}}$$

$$= \frac{240}{44.7213} = 5.49 > Z_{0.05} \text{ of } 1\%$$

H_0 is rejected. Die is biased.

- 2) In a sample of 1000 people in a state, 540 are rice eaters & the rest are wheat eaters. Can we at 1% & 5% L.O.S. that both rice & wheat are equally popular in the state.

$$P = \frac{1}{2}, n = 1000, q = \frac{1}{2}$$

for $x = 540$ (H_0 : rice & wheat are equally popular).

$$Z = \frac{540 - 1000 \cdot \frac{1}{2}}{\sqrt{1000 \cdot \frac{1}{2} \cdot \frac{1}{2}}} = -8.8$$

$$= \frac{40}{\sqrt{250}} = 40$$

$$Z = 2.5998$$

$Z > Z_{0.05}$ for 5% los. $\therefore H_0$ is rejected.

$$Z < Z_{0.05}$$
 for 5% los.

H_0 can be accepted.

- 3) A coin was tossed 400 times & heads turned up 225. Test the hypothesis that the coin is unbiased.

$$n = 400, x = 225, P = \frac{1}{2}, q = \frac{1}{2}$$

$$\sqrt{400 \cdot \frac{1}{2} \cdot \frac{1}{2}}$$

$$\sqrt{100} = 10$$

4. A manufacturer claims that only 4% of his products are defected. A random sample of 600 samples contains 36 defectives. Test the claim of the manufacturer at 5% level of significance.

$$n = 600, \quad n = 36, \quad p = 0.04, \quad q = 0.96$$

$$z = \frac{36 - 600(0.04)}{\sqrt{600 \times 0.04 \times 0.96}}$$

$$z = \frac{36 - 24}{\sqrt{24}} = \frac{12}{\sqrt{24}}$$

$$z = \frac{36 - 24}{4.8} = \frac{12}{4.8}$$

$$z = 2.5 > z_{0.05} \text{ for } 5\% \text{ L.O.S.}$$

The claim of manufacturer is rejected for 5% L.O.S. & accepted at 1% L.O.S.

5. A new medicine was given to 50 patients and found to be effective in 37 patients. Test the hypothesis that the medicine is effective in 80% patients.

$$n = 50, \quad n = 37, \quad p = 0.8, \quad q = 0.2$$

$$z = \frac{37 - 50(0.8)}{\sqrt{50 \times 0.8 \times 0.2}}$$

$$z = \frac{37 - 40}{\sqrt{8}}$$

PAGE NO:
DATE: / /

PAGE NO:
DATE: / /

$$z = 2.3, \quad z_{0.025} = 1.96$$

$$z = -1.0606$$

$$z \leq z_{0.05} \text{ for both } 5\%, 1\%, 10\% \text{ L.O.S.}$$

∴ The medicine is effective in 80% patients.

→ Test of significance for small samples ($n < 30$)

Degree of freedom:

It is the no. of independent values that a statistical analysis can estimate.
i.e. d.f. is the no. of values that are free to vary as we estimate parameters.

For Sample size less than 30, we do not assume that the sampling distribution is normal and the values given by the sample are sufficiently close to the population values.

→ Student's t-test

(Test of significance of the mean of a small sample)

Let x_i ($i = 1, 2, \dots, n$) be a random sample of size n from a normal population with mean μ and variance σ^2 .

PAGE NO.: / /
DATE: / /

Then $t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$ where $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

s^2 is an unbiased estimate of the population variance σ^2 and it follows t -distribution with $(n-1)$ d.f.

→ Questions of similar sets:

- Ques 1) A random sample of size 16 has a mean 53. The sum of squares of the deviations from mean is 135. Can the sample be regarded as taken from a population of mean 56? Also obtain 95% and 99% confidence limits for the mean of the population.

H_0 : population mean is 56, $\bar{x} = \mu$.

$$n = 16, \bar{x} = 53, \mu = 56, \sum (x_i - \bar{x})^2 = 135.$$

$$s^2 = \frac{1}{15} \times 135$$

$$s = \sqrt{\frac{1}{15} \sum (x_i - \bar{x})^2}$$

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{53 - 56}{\frac{3}{4}} = -4$$

$$t_c (5\%, 15 \text{ d.f.}) = 2.131$$

$$t_c (1\%, 15 \text{ d.f.}) = 2.947$$

PAGE NO.: / /
DATE: / /

$|t| = 4 > t_c$ (for 1% & 5% los.)

∴ H_0 is rejected for both 1% & 5%.
Hence i.e., $\bar{x} \neq \mu$.

→ Confidence Intervals

95% CI $\rightarrow -2.131 < t < 2.131$

$$-2.131 < \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} < 2.131$$

$$\frac{-s}{\sqrt{n}}(2.131) < \bar{x} - \mu < \frac{s}{\sqrt{n}}(2.131)$$

$$\bar{x} + \frac{2.131 s}{\sqrt{n}} > \mu > \bar{x} - \frac{2.131 s}{\sqrt{n}}$$

$$\mu \in (\bar{x} \pm 2.131 \frac{s}{\sqrt{n}})$$

$$\mu \in (53 \pm 2.131 \times \frac{3}{4})$$

$$\mu \in (51.9017, 54.5932)$$

99% CI \rightarrow

$$\mu \in (\bar{x} \pm 2.947 \frac{s}{\sqrt{n}})$$

$$\mu \in (53 \pm 2.947 \times \frac{3}{4})$$

$$\mu \in (50.4897, 55.2102)$$

2) A sample of 20 items has mean 42 units and s.d. 5 units. Test the hypothesis that it is a random sample from a normal population with mean 45 units.

$$H_0: \bar{x} = \mu$$

$$151.6 > 181.6 - 1.96 \times \frac{5}{\sqrt{20}}$$

$$n=20, \bar{x}=4, \mu=45, s=5 \Rightarrow s^2=25$$

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

$$S^2 = \frac{n-1}{n} S^2 \Rightarrow n-1 > (1.96)^2 \cdot 25$$

$$S^2 = \frac{20 \times 25}{19} = 26.3158$$

$$S = 5.0299$$

$$t = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{4 - 45}{5.0299/\sqrt{20}}$$

$$(-2.6153, 1.845)$$

$$t = -2.6153.$$

$t > t_c$ for 5% L.O.S.

$t \leq t_c$ for 1% L.O.S.

$\therefore H_0$ is rejected at 5%, los & accepted at 1% l.o.s.

$$\therefore 95\% \text{ CI} \rightarrow \mu \in \left(\bar{x} \pm 2.131 \frac{s}{\sqrt{n}} \right)$$

$$\mu \in \left(4 \pm 2.131 \frac{5.0299}{\sqrt{20}} \right)$$

$$\mu \in (4 \pm 2.444)$$

$$151.6 > 181.6 - 2.444 \times \frac{5.0299}{\sqrt{20}}$$

$$99\% \text{ CI} \rightarrow \mu \in \left(\bar{x} \pm 2.947 \frac{5.0299}{\sqrt{20}} \right)$$

$$151.6 > 181.6 - 2.947 \times \frac{5.0299}{\sqrt{20}}$$

$$\mu \in (0.6196, 7.3804)$$

3. 9 items of a sample have the values 45, 47, 50, 52, 48, 47, 49, 53, 51. Does it differ significantly from the assumed mean of 47.5.

$$H_0: \mu = \bar{x}$$

$$n=9, \mu=47.5, \bar{x}=49.111$$

$$S^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2 = \frac{1}{8} \times 54.8765$$

$$S^2 = 6.8596$$

$$S = 2.619$$

$$t = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{49.111 - 47.5}{2.619/\sqrt{9}} = 1.6111$$

$$t = 1.845 < t_c \text{ for } 5\% \text{ & } 1\% \text{ los}$$

$\therefore H_0$ is accepted for both 5% & 1% los

$$t_c(5\%, 8 \text{ d.f.}) = 2.306$$

$$t_c(1\%, 8 \text{ d.f.}) = 3.355$$

4. The heights of 10 boys in a class is 40, 67, 62, 68, 61, 68, 70, 64, 64, 66 inches. Is it reasonable to believe that the avg height of boys is 64 inches?

5. 10 students from a school have IQ 70, 120, 110, 101, 88, 83, 95, 98, 107, 100. Does this data support the assumption of avg (school) IQ to be 100?

$$n=10, \mu=100, \bar{x}=97.2$$

$$s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2 = \frac{1}{9} (\sum x_i^2 - n\bar{x}^2)$$

$$s = 14.2735$$

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{97.2 - 100}{14.2735/\sqrt{10}} = -2.18$$

$$t = 0.6203 < t_c(5\%, 9 \text{ d.f.})$$

$$t_c(5\%, 9 \text{ d.f.}) = 2.268$$

$$H_0: \bar{x} = \mu \quad (\text{is not accepted})$$

→ Two Sample t-test
For 2 samples of size n_1 & n_2

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{s^2_{xx} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$$s^2_{xx} = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1+n_2-2)}$$

$$s^2 = \frac{1}{n_1+n_2-2} \left\{ \sum (\bar{x}_i - \bar{x})^2 + \sum (y_j - \bar{y})^2 \right\}$$

$$H_0: \mu_1 = \mu_2$$

→ Questions:-

- i) Two samples give the foll. data

Sample A:-

44 44 56 46 47 47 58 53 49 55 46
57 51 60 53 52 50 54 56 58 59

Sample B:-

35 38 35 29 40 39 32 41 42 30 31
39

Test if they come from populations with the same mean.

$$\rightarrow H_0: \mu_1 = \mu_2 ; n_1 = 13, n_2 = 12$$

$$\bar{x} = 50.2307 \quad \bar{y} = 35.9166$$

$$s^2_{xx} = \frac{1}{n_1+n_2-2} \left\{ 306.307 + 226.915 \right\} = 23.17$$

$t \sim t(n_1 + n_2 - 2)$

$$t = \frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{\sqrt{s_x^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$$t = \frac{50.25 - 35.92}{\sqrt{23.17 \left(\frac{1}{13} + \frac{1}{12} \right)}} = 7.41$$

$t_c (5\%, n_1 + n_2 - 2 \text{ df})$

$$t = t_c (5\%, 23 \text{ df}) = 2.069.$$

$t > t_c \Rightarrow H_0 \text{ is rejected}$

2. Samples of 2 types of electric bulbs were tested for length of life & the following data was obtained.

Sample no. $n_1 = 8$ $n_2 = 7$ A sample
Mean $\bar{x}_1 = 123.4$ $\bar{y}_1 = 103.6$ B sample

SD $s_1 = 36$ $s_2 = 40$

From data, can you infer that type 1 is better than type 2?

$$\sum (x_i - \bar{x})^2 = n_1 s_1^2 = 8 \times 36^2 = 10368$$

$$\sum (y_i - \bar{y})^2 = n_2 s_2^2$$

$$= 7 \times 40^2 = 11200$$

$$s_x^2 = \frac{(8 \times 1) - 1}{8+7-2} \left\{ 10368 + 11200 \right\}$$

$$s_x^2 = 1659.077$$

$$t = \frac{198}{\sqrt{1659.077 \left(\frac{1}{8} + \frac{1}{7} \right)}}$$

$$t = 9.39$$

$$t_c (5\%, 13 \text{ df}) = 2.16$$

$t > t_c \therefore H_0 \text{ is rejected}$

3. A group of 10 children were fed on diet A & another group of 8 children were fed on diet B for a period of 6 months. The following increase in weights was recorded.

$$\begin{array}{cccccccccc} A: & 5 & 6 & 8 & 1 & 12 & 4 & 3 & 9 & 6 & 10 \\ B: & 2 & 3 & 6 & 8 & 10 & 1 & 2 & 8 \end{array}$$

Test whether diets A & B differ significantly regarding their effect on increase in weight.

$$\bar{x} = 6.4 \quad \bar{y} = 5 \quad n_1 = 10 \quad n_2 = 8$$

$$\sum (x_i - \bar{x})^2 = 110.24$$

$$\sum (y_i - \bar{y})^2 = 82$$

$$s_x^2 = \frac{1}{10+8-2} \left\{ 110.24 + 82 \right\}$$

$$s_x^2 = \frac{1}{16} \times 184.4 = 11.525$$

$$t = \frac{1.4}{\sqrt{11.525 \left(\frac{1}{10} + \frac{1}{8} \right)}} = \frac{1.4}{\sqrt{1.6103}} = 0.8694$$

$$t_c (5\%, 16 \text{ df}) = 2.12 \quad H_0 = H_a$$

$t < t_c \therefore H_0 \text{ is accepted}$.

A new process for producing synthetic diamond is viable only if the avg weight of the diamond > 0.5 carat. The weights of 6 diamond thus generated are $0.46, 0.51, 0.52, 0.48, 0.57, \& 0.54$ carat. Test the viability of the process.

$$\mu > 0.5, \bar{x} = 0.5133, n=6$$

$$\sum (x_i - \bar{x})^2 = 7.933 \times 10^{-3}$$

$$s^2 = \frac{1}{5} \times 7.933 \times 10^{-3} \Rightarrow s = \sqrt{1.5866 \times 10^{-3}}$$

$$t = 0.5133 - 0.5 = 0.8179.$$

$$t_{0.05} = 2.0398, \quad t_{0.05} = 2.0398, \quad t = 0.8179, \quad t < t_{0.05}, \quad H_0 \text{ is accepted}$$

$$H_0: \bar{x} = \mu, \quad H_a: \bar{x} \neq \mu \quad (\text{two-tailed test})$$

$$t(5\%, \text{los}, 5 \text{ df}) = \pm 2.015 \quad (\text{one-sided data})$$

$$t < t_{0.05}, \quad H_0 \text{ is accepted for } 5\% \& 1\% \text{ los.}$$

\rightarrow F-test.

It is used to calculate the ratio of variances of normally distributed statistics. Let $X = \{x_1, x_2, x_3, \dots, x_n\}$ & $Y = \{y_1, y_2, \dots, y_m\}$ be the values of 2 samples drawn from the same normal population. Then $F = \frac{s_1^2}{s_2^2}$

$$s_1^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \& \quad s_2^2 = \frac{1}{m-1} \sum_{i=1}^m (y_i - \bar{y})^2$$

: Larger value to be taken in the numerator (always)

\rightarrow F-test for comparing variance:-

1) Two samples of size 9 and 8 gives sum of squares of deviation from their resp. means as 160 inches² & 91 inches². Can these be regarded as drawn from the normal popl' with the same var?

$$S_1^2 = \frac{1}{8} (160) = 20, \quad S_2^2 = \frac{1}{7} (91) = 13$$

$$n=9, m=8, \sum (x_i - \bar{x})^2 = 160, \sum (y_i - \bar{y})^2 = 91$$

$$S_1^2 = \frac{1}{8} (160) = 20, \quad S_2^2 = \frac{1}{7} (91) = 13$$

$$F = \frac{20}{13} = 1.5384, \quad F_c(5, 1, 8, 7 \text{ df}) = 3.73$$

H_0 is accepted, as $F < F_c$. Thus the two popl' variances are equal.

a) Two independent samples have the foll' values

$$A: 28, 30, 32, 33, 33, 29, 34, \quad (7)$$

$$B: 29, 28, 27, 24, 23, 28, \quad (6)$$

Examine if they have been drawn from the normal popl' with the same variance.

$$\bar{x} = 31.2857, \bar{y} = 26.5, \sum (x_i - \bar{x})^2 = 31.4285, \sum (y_i - \bar{y})^2 = 29.5$$

$$s_1^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2 = \frac{1}{6} (31.4285) = 5.238$$

$$s_2^2 = \frac{1}{m-1} \sum (y_i - \bar{y})^2 = \frac{1}{5} (29.5) = 5.9$$

$$F = \frac{5.9}{5.238} = 1.1263, \quad F_c(5, 6 \text{ df}) = 4.39$$

$F < F_c \Rightarrow H_0$ accepted.

3) If 2 samples of size 10 & 14 have std dev 1.5 & 1.2, can we say that they are drawn from normal popl with the same variance?

$$\rightarrow n_1 = 10, n_2 = 14, S_1 = 1.5, S_2 = 1.2$$

$$S_p^2 = \frac{n_1 S_1^2 + n_2 S_2^2}{n_1 + n_2 - 2} = \frac{10 \times (1.5)^2 + 14 \times (1.2)^2}{10 + 14 - 2} = 2.5$$

$$S_p^2 = \frac{n_1 S_1^2}{n_2 - 1} = \frac{10 \times (1.5)^2}{14 - 1} = 1.5507$$

$$F = \frac{2.5}{1.5507} = 1.6121 \quad F_c(5\%, 9, 13) = 2.7$$

$F < F_c \Rightarrow H_0$ accepted!

4) A teacher in 2 classes A and B, the same subject, class A has 16 students & 25 students. In an exam, although there was no big diff in mean grades, class A has a std dev of 9 & B has 12. Can we conclude that variability of B > A.

→ Chi Square test

It is used for

- Independence of attributes
- goodness of fit

* Most of the tests were based on the assumption that the sample was drawn from a normal popl. But there are situations when the data doesn't follow normal dist and the observations are recorded only as presence or absence of an attribute. In such a situation, we use non-parametric tests of significance.

$$X_i^2 = (O_i - E_i)^2$$

$O_i \rightarrow$ Observed frequency

$E_i \rightarrow$ Expected "

$$\chi^2 = \sum_i X_i^2$$

when null hypothesis is true, χ^2 follows chi-square dist with $(r-1)(c-1)$ d.f.

where, r = no. of rows

c = no. of columns

→ Questions:

- In a health survey, 400 individuals were asked if they were vaccinated against flu and if they subsequently had ~~exact~~ flu. The full data was

obtained. Test if vaccination prevents flu.

	Vaccinated		Not vax	
got flu	60	85	145	
Not got flu	190	65	255	
	250	150	400	

H_0 : vaccination & catching flu are independent.

Attribute	O_i	E_i	$\chi^2 = \frac{(O_i - E_i)^2}{E_i}$
(vacc & got flu)	60	$145 \times 250 = 90.625$	10.347
		400	
(Not vacc & got flu)	85	$150 \times 145 = 54.375$	17.248
		200	
(vacc & not got flu)	190	$250 \times 255 = 159.375$	5.8848
		400	
(Not vacc & not got flu)	65	$150 \times 255 = 95.625$	9.808
		400	
			43.29

$$\chi^2 = \sum_i \chi^2_i = 43.29.$$

$$df = (2 \times 1) (2 - 1) = 1$$

$$\chi^2 (1 d.f.) = 3.84.$$

0.95

H_0 is rejected.

- Q) Given the data below, perform χ^2 test to check if hair colour and eye colour are dependent.

Hair colour:

Eye colour	Black		Grey	Total
	Black	Blue		
Black	40	20	60	
Blue	20	30	50	
Brown	60	30	90	
	120	80	200	

H_0 : Hair colour & Eye colour are independent.

Attribute	O_i	E_i	$\chi^2_i = \frac{(O_i - E_i)^2}{E_i}$
Black/Black	40	$120 \times 60 = 36$	0.4444
grey/Black	20	200	
grey/Black	20	$60 \times 80 = 24$	0.6667
Black/Blue	20	$50 \times 120 = 30$	3.3333
grey/Blue	30	$50 \times 80 = 20$	5.0
Black/Brown	60	$120 \times 90 = 54$	0.6667
grey/Brown	30	$80 \times 90 = 36$	1.0

11.1111.

$$\chi^2 = \sum_i \chi^2_i = 11.1111 > 3.84$$

$$\chi^2_{(0.95)} \{ (3-1)(2-1) df \} = \chi^2_{(0.95)} (2 df) = 5.99$$

$\therefore H_0$ is rejected.

- a) The no. of demands for a particular spare part in a shop is found to vary day by day. In a sample study, the following data is obtained.

Day	Mon	Tue	Wed	Thu	Fri	Sat
No. of parts demanded	124	125	110	120	126	115

$$d.f = (r-1)(c-1) = (2-1)(2-1) = 1$$

$$\chi^2_c (0.95, 1 d.f) \approx 5.01$$

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} = 1.6832$$

Attribute	O _i	E _i	χ^2_i
1	124	120	0.133
2	125	120	0.2083
3	110	120	0.8333
4	120	120	0
5	126	120	0.3
6	115	120	0.2083
			1.6832

$$\chi^2 < \chi^2_{0.95, 1 d.f}$$

H_0 is accepted.

- (Q) A sample analysis of exam results of 500 students was made. Is as follows. 200 failed, 170 had III class, 90 had II class, 20 had I class. Do the following the ratio 4:3:2:1?

Total = 500, 4:3:2:1

	O _i	E _i	χ^2_i
Fail	200	200	
III	170	150	
II	90	150	
I	20	50	

$$\chi^2_c (0.95, 3 d.f) = 7.81$$

- (Q) A survey of 800 families with 12 children each regarde the foll data

No. of boys	0	1	2	3	4	5
No. of girls	32	178	290	236	64	8
Families						

Is the result inconsistent with the hypothesis that a girl or a boy is equally likely to be born?

$$p=1/2, q=1/2$$

B	G	O _i	E _i	χ^2_i
0	4	32	$4C_0 (1/2)^0 (1/2)^4 = 1/16 \times 800 = 50$	
1	3	178	$4C_1 (1/2)^1 (1/2)^3 = 4/16 \times 800 = 200$	
2	2	290	$4C_2 (1/2)^2 (1/2)^2 = 6/16 \times 800 = 300$	
3	1	236	$4C_3 (1/2)^3 (1/2)^1 = 4/16 \times 800 = 200$	
4	0	64	$4C_4 (1/2)^4 (1/2)^0 = 1/16 \times 800 = 50$	

$$\chi^2_c : 6.48 \quad 2.42 \quad 0.3333 \quad 0.48 \quad 3.92$$

$$\sum O_i^2 = \chi^2 = 19.633$$

$$\chi^2_c (0.95, (5-1)(2-1) d.f) = \chi^2_c (0.95, 4 d.f) = 9.49$$

H_0 is rejected.

Q) A die is thrown 276 times and the results of these throws are shown below.

No. appeared	Actual Freq	Expt Freq	χ^2_i
1	44	46	0.0869
2	52	46	0.7826
3	42	46	0.3478
4	50	46	0.3478
5	49	46	0.1956
6	39	46	1.0652

$$\chi^2_c(0.95, 5 \text{ df}) = 11.1$$

$$\sum x_i^2 < \chi^2_c \quad \text{H}_0 \text{ is accepted}$$

Q. Test the given data for goodness of fit for Poisson dist

x	0	1	2	3	4
f _i	122	60.88	15	2.81	8.62

$$\mu = \frac{\sum x_i f_i}{\sum f_i} = \frac{100}{200} = 0.5$$

$$P(x) = \frac{e^{-\mu} \mu^x}{x!}$$

$$P(x=0) = e^{-0.5} \frac{(0.5)^0}{0!} = 0.6065$$

$$P(x=1) = e^{-0.5} \frac{(0.5)^1}{1!} = 0.3033$$

$$P(x=2) = e^{-0.5} \frac{(0.5)^2}{2!} = 0.0758$$

$$P(x=3) = e^{-0.5} \frac{(0.5)^3}{3!} = 0.0126$$

$$P(x=4) = e^{-0.5} \frac{(0.5)^4}{4!} = 0.00158$$

~~$$P(x>5) = e^{-0.5} \frac{(0.5)^5}{5!} = 0.000258$$~~

$$E(O) = 0.60 = 0.00001 \quad 0.00001$$

$$x_0 = 0.00001 \quad E(O) = X^2_{0.00001}$$

$$0 = 122 - 121.301 = 0.0044$$

$$1 = 60 - 60.66 = 0.0072$$

$$2 = 15 - 15.16 = 0.0016$$

$$3 = 2 - 2.52 = 0.1073$$

$$4 = 1 - 3.16 = 1.4805$$

$$\sum x_i^2 < \chi^2_c \quad \text{H}_0 \text{ is accepted}$$

We lose 1 d.f. in calculating μ .

$$\therefore \text{d.f.} = (5-1)(2-1) - 1 = 3 \text{ d.f.}$$

$$\chi^2_c(0.95, 3 \text{ d.f.}) = 7.81$$

$$\sum x_i^2 < \chi^2_c$$

$\therefore \text{H}_0$ is accepted.

Q) The table gives the no. of good and bad parts produced by each of the 3 shifts in a factory. Test if the prodⁿ of bad parts is independent of the shifts.

shifts	Good	Bad.	
1	82960	40	1000
2	940	50	990
3	950	60	1100
	2850	150	3000

shifts	Bad. O _i	E _i	X _i ²
1	960	$\frac{1000 \times 2850}{3000} = 960$	0.1053
2	940	$\frac{990 \times 2850}{3000} = 940.5$	0
3	950	$\frac{1100 \times 2850}{3000} = 1045$	0

$$\chi^2_{c} (0.95, 2 \text{ df}) = 5.99$$

$$2X^2 < X^2_{c(0.95)}$$

$\therefore H_0$ is accepted.

Q) Fit a Poisson distribution for the following data and test the goodness of fit.

x	0	1	2	3	4	5	6	Total
f	273	70	30	7	7	2	1	390

$$\mu = \frac{\sum f x}{\sum f} = \frac{195}{390} = 0.5$$

$$P(x=0) = \frac{e^{-0.5} (0.5)^0}{0!} = 0.6065 \times 390 = 236.535$$

$$P(x=1) = \frac{e^{-0.5} (0.5)^1}{1!} \times 390 = 118.287$$

$$P(x=2) = 49.562$$

$$P(x=3) = 4.914$$

$$P(x=4) = 0.6162$$

$$P(x=5) = 0.06162$$

$$P(x=6) = 0.00513$$

x	O _i	E _i	X _i ²
0	273	236.535	5.6065
1	70	118.287	19.7116
2	30	29.562	0.0065
3	7	4.914	0.8851
4	1	0.6162	66.1358
5	0	0.06162	60.9756
6	0	0.00513	192.9369

$$\chi^2_c (0.95, 5 \text{ df}) = 11.8$$

$\sum X_i^2 > 11.8$
 $\therefore H_0$ is rejected.

	Face appeared on dice.	O _i	E _i	X _i ²
1		85	$\frac{1}{6} \times 300 = 50$	0.5
2 or 3		108	$\frac{2}{6} \times 300 = 100$	0.64
4 or 5		94	$\frac{2}{6} \times 300 = 100$	0.36
6		43	$\frac{1}{6} \times 300 = 50$	0.98
		300		2.48

$$\chi^2_c (0.95, 3 \text{ df}) = 7.81$$

$X^2_i < \chi^2_c$ $\therefore H_0$ is accepted.

Q.	No. of heads	O _i	E _i at P ₀	χ^2	
				$\sum C_0(Y_2)^5 (Y_2)^0 = 10$	2.5
	0	15		$\sum C_0(Y_2)^5 (Y_2)^1 = 50$	0.5
	1	45		$\sum C_0(Y_2)^5 (Y_2)^2 = 100$	2.25
	2	85		$\sum C_0(Y_2)^5 (Y_2)^3 = 100$	0.25
	3	95		$\sum C_0(Y_2)^5 (Y_2)^4 = 50$	2
	4	60		$\sum C_0(Y_2)^5 (Y_2)^5 = 10$	10
	5	20			
		320			17.5

$$\chi^2_{\text{cal}}(0.95, \text{sd}) = 11.1$$

$\chi^2_{\text{cal}} > \chi^2_{\text{tab}}$ H_0 is rejected.

Q. Attribute O_i E_i S.E.P. χ^2

Red $134.8 - 0.9/16 \times 64 = 367.13$ 0.1111

Black $310.1 - 3/16 \times 64 = 122.13$ 0.3333

white $20.2 - 0.4/16 \times 64 = 116.13$ 1

100P. S.P. $\chi^2_{\text{cal}} = 1.4444$

$$\chi^2_{\text{cal}}(0.95, 2 \text{ df}) = 5.99 \text{ P.} \rightarrow$$

H_0 is accepted.

Ans.

$$2.0 \quad 0.2 = 0.05 \times 4 \quad 2.8$$

$$1.9 \quad 0.01 = 0.05 \times 2 \quad 0.1$$

$$1.8 \quad 0.01 = 0.05 \times 3 \quad 0.15$$

$$1.7 \quad 0.2 = 0.05 \times 3.1 \quad 0.2$$

$$1.6 \quad 0.02 = 0.05 \times 3.2 \quad 0.2$$