

an example of mean standard deviation

Two sample Z-test

- This test is used to check whether the two samples taken from the same population or not.

→ two samples are same or not

Two sample Z-test assumptions

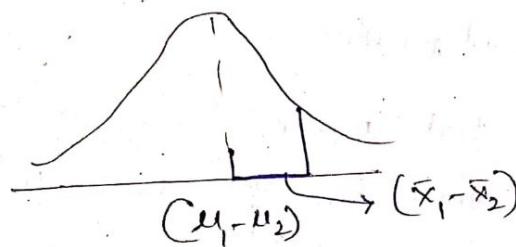
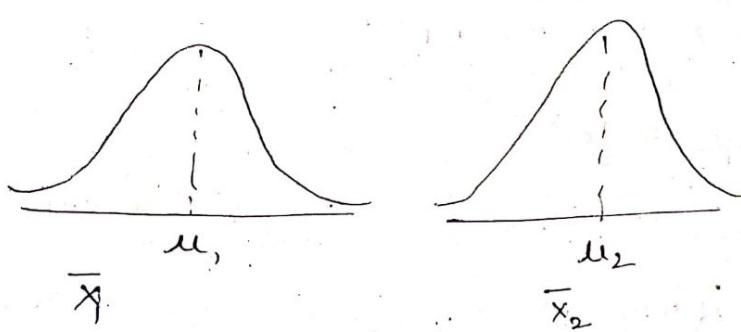
2 goals two sample Z-test:

- A hypothesis test that is used to compare two sample groups to determine if they have originated from same population or from different populations.
- Two sample Z-test requires the standard deviation to be known & the sample size to be larger than 30 and should population should follow normal distribution.

Two Sample Z-test assumptions:

- The sample size (say n_1 & n_2) of two samples drawn from two populations are large (say at least 30) and corresponding standard deviations σ_1 , σ_2 are known.
- Samples are drawn from 2 normally distributed populations with standard deviations with σ_1 , σ_2 known.
- Assume that μ_1 & μ_2 are population means. Our interest is to check a hypothesis on difference $\mu_1 - \mu_2$.

- If \bar{x}_1 & \bar{x}_2 are estimated means from two samples drawn from 2 populations. The statistic $(\bar{x}_1 - \bar{x}_2)$ follows standard normal distribution with mean $(\mu_1 - \mu_2)$ and standard deviation $\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$.



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

Ques The marketing specialization students earns 5000 Rupees more per month than operation management students atleast.

Specialization	Sample size	salary	standard deviation
Marketing	120	67500	7200
Operations	45	58950	4600

$$\bar{x}_1 = 67500$$

$$\bar{x}_2 = 58950$$

$$n_1 = 120$$

$$n_2 = 45$$

$$H_0: \mu_1 -$$

$$H_1: \mu_1 -$$

$$Z =$$

$$\Rightarrow$$

$$\bar{x}_1 = 63500 \quad \bar{x}_2 = 58950$$

$$\sigma_1 = 1440$$

$$n_1 = 120$$

$$s_1 = 3500$$

$$s_2 = 4,600$$

$$H_0: \mu_1 - \mu_2 \leq 5000$$

$$H_1: \mu_1 - \mu_2 > 5000$$

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

$$\therefore \mu_1 - \mu_2 = 5000$$

$$\Rightarrow \frac{(63500 - 58950) - (5000)}{\sqrt{\frac{(7200)^2}{120} + \frac{(4,600)^2}{45}}}$$

$$\sqrt{\frac{(7200)^2}{120} + \frac{(4,600)^2}{45}}$$

$$\Rightarrow \frac{0.8550 - 5000}{\sqrt{\frac{518400}{120} + \frac{211600}{45}}}$$

$$\sqrt{\frac{518400}{120} + \frac{211600}{45}}$$

$$\Rightarrow \frac{3.550}{4320 + }$$

$$= \underline{\underline{3.737}}$$

$$\begin{array}{r} 72 \times 72 \\ \hline 144 \\ 504 \\ \hline 5484 \\ \hline 23 \\ 46 \times 46 \\ \hline 276 \\ 184 \\ \hline 2116 \end{array}$$

$$\begin{array}{r} 12) 51840 (432 \\ \hline 48 \\ 38 \\ 36 \\ \hline 24 \\ 24 \\ \hline 0 \end{array}$$

$$\begin{array}{r} 45) 211600 (464 \\ \hline 180 \\ 316 \\ 315 \\ \hline 100 \\ 90 \\ \hline 100 \end{array}$$

Tying

Q10) Typing speed on a pc

Men & women typing speeds are different.

	men	women
\bar{x}	65 wpm	68 wpm
s	10 wpm	14 wpm
n	50	60

consider $\alpha = 1\%$.

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$\Rightarrow \frac{(\bar{x}_1 - \bar{x}_2) - (0)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$\Rightarrow \frac{(65 - 68)}{\sqrt{\frac{(10)^2}{50} + \frac{(14)^2}{60}}} \Rightarrow \frac{-3}{\sqrt{2 + 3.26}} \Rightarrow \frac{-3}{\sqrt{5.26}} \\ \Rightarrow \frac{-3}{2.29} \Rightarrow \underline{\underline{-1.30}}$$

$$\begin{aligned} &\frac{3}{19.6} \\ &\frac{7}{66.3015} \\ &15) \frac{49(3.26)}{40} \\ &\frac{30}{100} \end{aligned}$$

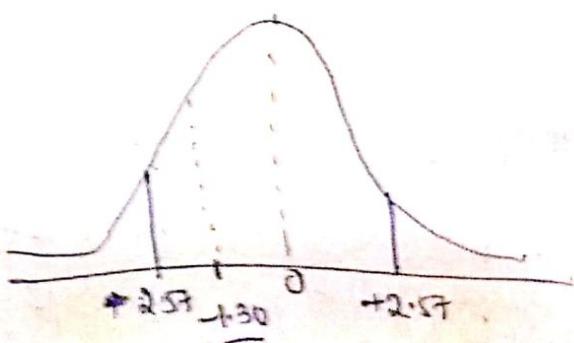
confidence bel 99%.

→ two tail test Z -score = ± 2.57

-1.30 is in accepted Regim.

so, H_0 will be accepted

$$\text{i.e., } \mu_1 = \mu_2$$



Ques Are the machine tools manufactured by X & Y different
w.r.t how long they last.

<u>sample</u>	Company X	Company Y
X	16.2 weeks	15.9 weeks
Y	0.2 weeks	0.2 weeks
Z	40	40

$$d = 0.08 \Rightarrow 3\%$$

SOL

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (u_1 - u_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$\sqrt{\frac{(0.2)^2}{40} + \frac{(0.2)^2}{40}}$$

$$\Rightarrow \frac{(16.2 - 15.9)^2 - (0)}{\sqrt{\frac{(0.2)^2}{40} + \frac{(0.2)^2}{40}}}$$

$$\sqrt{\frac{(0.03)^2}{40} + \frac{(0.03)^2}{40}}$$

$$\Rightarrow \frac{(0.03)^2}{\sqrt{\frac{0.04 + 0.04}{40}}} = \frac{(0.03)^2}{\sqrt{\frac{0.08}{40}}} = \frac{(0.03)^2}{\sqrt{0.002}}$$

$$\Rightarrow \frac{0.09}{\sqrt{\left(\frac{0.08}{40}\right)}} \Rightarrow \frac{0.09}{\sqrt{\left(\frac{8}{4000}\right)}} = \frac{0.09}{\sqrt{\frac{1}{500}}} \Rightarrow 0.09 \times \sqrt{500} = 0.09 \times 22.36 = 2.01$$

\Rightarrow

P20 Who earns more married (H0) uncorrelated

	married	unmarried	
X	\$ 639.60	\$ 658.20	Condition $\sigma = 66.04$
S	60	70	
N	40	60	

Sol H_0 : earnings of married & unmarried are same.

H_1 : earnings of married & unmarried are not same.

Z-statistic

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sigma / \sqrt{n}} \Rightarrow \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$\frac{65}{400} - \frac{150}{200} \Rightarrow 1.24$$

$$\Rightarrow \frac{(639.60 - 658.20) - 0}{\sqrt{60^2/40 + 70^2/60}}$$

125

H_0 woman lives longer married & unmarried

	married	unmarried	$\neq 7.04 \text{ years}$	
\bar{x}	78.5 years			$d = 0.01$
s	14.0	16.0		
n	140	160		

H_0

H_1

F-statistic

$$= \frac{(x_1 - x_2)^2 - (u_1 - u_2)^2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$\Rightarrow \frac{(78.5 - 77.0)^2 - (0)^2}{\sqrt{\frac{(14)^2}{140} + \frac{(16)^2}{160}}}$$

$$\Rightarrow (1.5)$$

$$\sqrt{\frac{196}{140} + \frac{256}{160}}$$

$$\Rightarrow \frac{1.5}{\sqrt{3}}$$

$$\Rightarrow$$

$$\frac{3.9 \times 14 + 16}{140} = \frac{16}{10}$$

$$= 1.4 + 1.6$$

$$\Rightarrow \frac{1.4 + 1.6}{2} = \frac{3.0}{2} = 1.5$$

$$\frac{2 \cdot 3}{2} = 1.5$$

$$\left(\frac{3}{2}\right)^2 = \frac{9}{4} = 2.25$$

Statis

Paired Sample t-test

- applied on a single sample, but at different conditions

(i) two different situations

Ex growth of a student before taking drug and after taking drug

→ In many cases we would like to analyse whether an intervention (or treatment) such as training programs, treatment

for specific illness may have significantly change population

Parameter. Such as mean (or) proportion, before & after intervention

- In pair t-test, the data related to parameter is captured twice. Once before intervention and once after intervention

- Ex: 1) Body weight before & after attending Yoga training
2) cholesterol level before & after attending meditation training.

...

Assume that the mean difference in the estimated parameter before & after treatment is ' d ' and the corresponding standard deviation difference is S_d . Let μ_d be the hypothesised mean difference. Then statistic 't'

$$t = \frac{d - \mu_d}{S_d / \sqrt{n}}$$

Assume d follows normal distribution then statistic 't' follows 't' distribution with $(n-1)$ degrees of freedom

Q10 A Researcher believes that people drink more coffee on Monday than other days of week. Based on sample of 50 coffee drinkers, the mean difference was estimated as 14 ml and corresponding standard deviation 8.5 ml. Conduct a hypothesis test at $\alpha = 0.1$ to claim that people drink on average 10 ml more coffee on Monday compared to other days of week.

$$\text{Sol} \quad n = 50$$

$$D = 14 \text{ ml}$$

$$S_d = 8.5 \text{ ml}$$

$$M_D = 10 \text{ ml}$$

$$\alpha = 0.1$$

$$H_0: M_D \leq 10$$

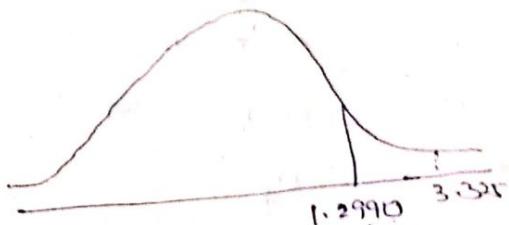
$$H_1: M_D > 10$$

$$t = \frac{D - M_D}{S_d / \sqrt{n}} = \frac{14 - 10}{8.5 / \sqrt{50}} = 3.3275$$

$$H_1: M_D > 10$$

$$H_0: M_D \leq 10$$

$$H_1: M_D > 10$$



H_0 : it has to be rejected.

so, people will have more coffee on Monday than other days.

P30 The difference in avg. weekly consumption for students before & after breakup is 11.5 and the corresponding standard deviation difference is 95.67 for 20 candidates.

conduct a hypothesis test to check whether weekly consumption is more after breakup i.e., ($H_1: \mu_d > 0$) at 95% confidence

Soln

$$t_d \approx 11.5$$

$$D = 11.5$$

$$S_d = 95.67$$

$$n = 20$$

$$\alpha = 0.05\%$$

$$\text{confidence level} = 95\%$$

Eqn

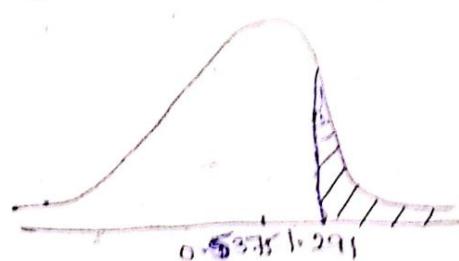
$$t = \frac{D - \mu_d}{S_d / \sqrt{n}} \Rightarrow \frac{11.5 - 0}{95.67 / \sqrt{20}} = 0.5375$$

t critical value at degree of freedom - 19

$$\alpha = 0.05$$

so, 0.5375 is in acceptance region

so, H_0 is accepted i.e., there is no difference.



Two sample t-test

(a) Differencing Two population means when population standard deviations are unknown, and believed to be equal.

If we want to estimate difference in two population means when standard deviation of the population are unknown then we need to estimate from the samples drawn from two populations.

An additional assumption is that the standard deviation of two populations are equal (otherwise unknown). Then sampling distribution of the difference in estimated means ($\bar{x}_1 - \bar{x}_2$) follows t-distribution with

degree of freedom $(n_1 + n_2 - 2)$ with mean $(\mu_1 - \mu_2)$ with standard deviation $\sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$.

here s_p^2 is the pooled variance.

s_p^2 = Pooled variance

s_p = Pooled standard deviation

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{(n_1 + n_2 - 2)} \quad (\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)$$

the corresponding t-statistic is $t = \frac{\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)}{\sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$

Ques A company makes a claim that children who drink their health drink will grow taller than the children who do not drink that health drink.

Group	Sample size	Increase in height	standard deviation
drink health drink	80	1.6 cm	1.1 cm
do not drink health drink	80	6.3 cm	1.3 cm

at $\alpha = 0.05$, test whether the increase in height of the children who drink the health drink is atleast 1.2 cm

Sol given. $n_1 = 80 \quad n_2 = 80$
 $\bar{x}_1 = 7.6 \quad \bar{x}_2 = 6.3$
 $s_1 = 1.1 \quad s_2 = 1.3$

$$H_0 : \Delta u_0 (u_1 - u_2) \leq 1.2$$

$$H_1 : (u_1 - u_2) > 1.2$$

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

$$= \frac{(1.3) - 1.2}{\sqrt{1.45 \left(\frac{1}{80} + \frac{1}{80} \right)}}$$

$$\Rightarrow \frac{0.1}{\sqrt{1.45 \times \frac{1}{40}}} \Rightarrow \underline{0.5252}$$

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

$$\Rightarrow \frac{79 \cdot (1.1)^2 + 79 \cdot (1.3)^2}{160-1}$$

$$\Rightarrow \frac{79 \cdot (1.21) + 79 \cdot (1.69)}{159}$$

$$\underline{1.45}$$

$$t\text{-statistic} = 0.5252$$

t-critical value for t-test test,

where $\alpha = 0.05$ & degree of freedom

$$(n_1 + n_2 - 2) = 158$$

$$\text{is } 1.6546$$

critical value of t-test is ≈ 1.6546

As t-statistic is not in the rejection region.

So, accept H_0 .

i.e., the difference is ≤ 1.6546 , accepting H_0 .

Two sample t-test is not applicable in this case.

~~t-test~~ (b) Difference in two population means when population

standard deviations are unknown and not equal:

two sample t-test with unequal variance

If we want to estimate difference in two population

means when standard deviations of the two populations

are unknown and unequal.

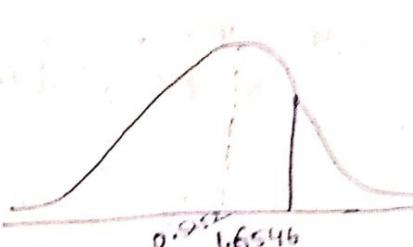
We need to estimate standard deviation from the samples drawn from these two populations.

Then the sampling distribution of the difference in estimated

means $(\bar{x}_1 - \bar{x}_2)$ follows a t-distribution with mean $(\mu_1 - \mu_2)$

and standard deviation $S_u = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$, the corresponding

degree of freedom is given by $df = \left\lfloor \frac{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}{\frac{(s_1^2/n_1)^2}{n_1-1} + \frac{(s_2^2/n_2)^2}{n_2-1}} \right\rfloor$



- Mean t -

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (u_1 - u_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

~~Ques~~ A Researcher is interested in finding the avg. duration of marriage based on educational qualifications of couples. Two groups were considered for study. Group 1 consisted of couple with no bachelors degree (both partners) and Group 2 consists of couple who both have bachelors degree or higher.

Group	Sample size	duration of marriage	standard deviation of sample
No degree	120	10.1 years	2.4 years
with degree	100	9.5 years	3.1 years

at $\alpha = 0.05$ test whether the avg. duration of marriage is same for couples with no bachelors degree as compared to couples with bachelors degree.

Sol $n_1 = 120 \quad n_2 = 100$

$$\bar{x}_1 = 10.1 \text{ yrs} \quad \bar{x}_2 = 9.5 \text{ yrs}$$

$$s_1 = 2.4 \text{ yrs} \quad s_2 = 3.1 \text{ yrs}$$

$$u_1 - u_2 \leq 0$$

$$\text{H}_0: u_1 - u_2 = 0$$

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (u_1 - u_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$\Rightarrow (10.1 - 9.5) = 0$$
$$\sqrt{\frac{(2.4)^2}{120} + \frac{(3.1)^2}{100}}$$

$$\frac{2.4 \times 2.4}{96}$$

$$\frac{H^2}{5.76}$$

$$\Rightarrow 0.6$$
$$\sqrt{\frac{5.76}{120} + \frac{9.61}{100}}$$

$$\frac{3.1 \times 3.1}{96}$$

$$\frac{9.3}{96}$$

$$\Rightarrow \sqrt{\frac{0.48}{10} + \frac{9.61}{100}}$$

$$\frac{12.25.76}{48} \frac{64.8 \times 12}{5.76}$$

$$\approx 1.5805$$

$$d_f = S_u = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$d_f = \left[\frac{(s_1)^4}{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2} \right] = \frac{0.0207}{0.000113}$$

$$= [\frac{33}{184.33}] \frac{184}{184} \approx 184$$

$$\alpha = 0.05 \quad \& \quad d_f = [184.33] \approx 184$$

The critical value of t for $\alpha = 0.05$ & degree of freedom (d_f)

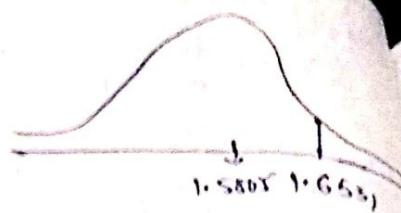
= 184. is 1.6531

t-critical value = 1.6531

t-statistic value = 1.5805

accept it i.e., H_0 is accepted.

i.e., $\mu_1 - \mu_2 \leq 0$.



Since the t-statistic is less than critical value of 't', we retain null hypothesis. i.e., the difference in duration of marriage between two groups is less than & equal to zero.

11/9/19

Non parametric test

good fit test

Parametric tests

Z-tests } mean

t-tests

F-test — proportion

Day	M	T	W	Th	F	S
-----	---	---	---	----	---	---

Expected %	10	10	15	20	30	15
------------	----	----	----	----	----	----

Observed	30	14	34	45	57	20
----------	----	----	----	----	----	----

Expected	20	20	30	40	60	30
----------	----	----	----	----	----	----

value

(Expected-Observed) 10 -6 4 5 -3 -10
(E) (O)

$(E - O)^2$ 10^2 $(-6)^2$ 4^2 5^2 $(-3)^2$ $(-10)^2$

Non-parametric test

In previous sections we discussed Z -test, t -test are called parametric tests. since the objective is to infer about population parameters such as mean and population parameters.

To conduct Z -test, t -test we need summary statistics (mean & standard deviation) not necessarily the entire distribution.

In parametric test, we assume that the population follows Normal distribution, Non-parametric tests are distribution free tests. since they do not have assumptions about distribution population.

The major difference between parametric & Non-parametric tests.

- parametric tests needs only values of parameters and knowledge about distribution
- In non-parametric test we need entire distribution of the data.
- Importantly data may not follow any parametric distribution such as Normal distribution
- Also test is not about the population parameter, but about characterizing the entire distribution.
- Non-test parameters

When do population non-parametric test is used

- The test is not about the population parameter such as Mean & standard deviation
- The method / test doesn't require assumption about population distribution.

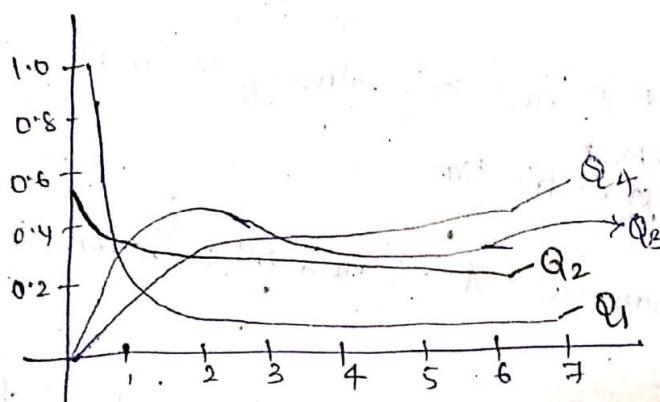
χ^2 -distribution chi-square (χ^2)

- let us take a random variable 'x' if it is coming from standard Normal distribution

$$X \sim N(0,1)$$

Q_1 is another random variable

$Q = x^2$, the distribution of Q , random variable is ~~χ^2~~ χ^2 -distribution with 'degree of freedom'.



- If we take 2-samples x_1, x_2, \dots, x_n ,

$Q_2 = x_1^2 + x_2^2$, which will follows χ^2 -distribution with $df = 2$

$$Q_3 = x_1^2 + x_2^2 + x_3^2$$

follows χ^2 -distribution with $df = 3$

A Restaurant owner believes that his restaurant will incoming customers expected distribution would be

day	M	T	W	Th	F	S
Expected %.	10	10	15	20	30	15

But on one week observed values are

Observed values	30	14	34	45	53	20
						respectively

At $\alpha = 0.05$, check whether owners expected customer's distribution correct or not, with appropriate hypothesis test.

Sol H_0 : owner's distribution is correct

H_1 : owner's distribution is not correct

χ^2 - statistic

Expected %.	10	10	15	20	30	15
observed	30	14	34	45	53	20

$$\frac{10 \times 200}{100}$$

Expected value	20	20	30	40	60	30

$$\chi^2\text{-statistic} = \frac{(30-20)^2}{20} + \frac{(14-20)^2}{20} + \frac{(34-30)^2}{30} + \frac{(45-40)^2}{40} + \frac{(53-60)^2}{60} + \frac{(20-30)^2}{30}$$

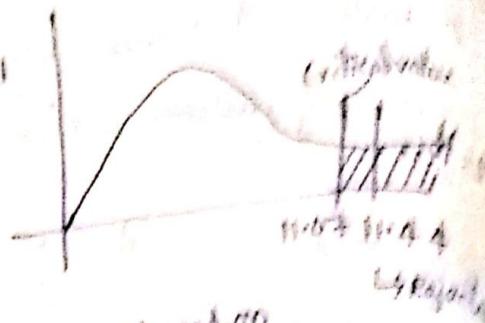
$$\Rightarrow \frac{10^2}{20} + \frac{6^2}{20} + \frac{4^2}{30} + \frac{5^2}{40} + \frac{3^2}{60} + \frac{10^2}{30}$$

$$\Rightarrow \frac{5}{2} + \frac{9}{5} + \frac{16}{15} + \frac{25}{8} + \frac{9}{60} + \frac{100}{30} = 11.44$$

χ^2 -distribution with $\alpha=0.05$ is 11.07

$\because \chi^2$ -statistic value $> \chi^2$ -critical value

We reject null hypothesis.



→ owner's distribution is not a good fit based on significance level 0.05

Ques 13

Hanuman Airlines operated Indian airlines to several Indian cities. One of the problems Hanuman Airlines faces is that preferences given by the passengers. Captain Cook the operations manager of Hanuman Airlines believes that 35% of passengers prefer vegetarian food, 40% passengers prefer non-vegetarian food, 20% low calorie food & 5% requested for diabetic food.

A sample of 500 passengers were chosen to analyse the food preferences

Food Type	Veg	Non-Veg	Low Calorie	Diabetic
No of passengers	190	185	90	35

Conduct a χ^2 test to check whether Captain Cook's belief is true at $\alpha = 0.05$.

Sol H_0 : Captain Cook's belief is true

H_1 : Captain Cook's belief is not true

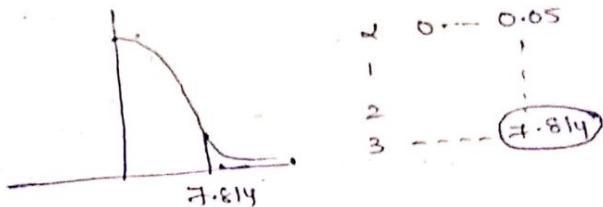
~~X²~~-statistic

Food type	Veg	non-veg	lowcalorie	diabetic	$\epsilon E = 500$
observed	170	185	90	35	
Expected %	35	45	20	5	
Expected value	175	225	100	25	

degrees

degree of freedom = 3

χ^2 -critical value with $df = 3$, & at $\alpha = 0.05$ is 7.814.



χ^2 -statistic $\left(\frac{O-E}{E}\right)$

$$\Rightarrow \left(\frac{190-175}{175} \right)^2 +$$

$$\left(\frac{(225-185)^2}{225} \right) + \left(\frac{(100-90)^2}{100} \right) + \left(\frac{(35-25)^2}{25} \right)$$

$$\Rightarrow \frac{225}{175} + \frac{(40)^2}{225} + \frac{(10)^2}{100} + \frac{(10)^2}{25}$$

$$\left(\frac{45}{18} \right)^2 = 9$$

$$\Rightarrow 1.2857 + 7.11 + 1 + 4$$

$$\Rightarrow \underline{15.3967}$$

$$7.9 (1.2857) \\ \frac{7}{20} \\ \frac{14}{60} \\ \frac{56}{40} \\ \frac{35}{50}$$