

Chi-Square Test

The chi-square test is an important test amongst the several tests of significance developed by statisticians. Chi-square, symbolically written as χ^2 (Pronounced as Ki-square), is a statistical measure used in the context of sampling analysis for comparing a variance to a theoretical variance. As a non-parametric* test, it “can be used to determine if categorical data shows dependency or the two classifications are independent. It can also be used to make comparisons between theoretical populations and actual data when categories are used.”¹ Thus, the chi-square test is applicable in large number of problems. The test is, in fact, a technique through the use of which it is possible for all researchers to (i) test the goodness of fit; (ii) test the significance of association between two attributes, and (iii) test the homogeneity or the significance of population variance.

CHI-SQUARE AS A TEST FOR COMPARING VARIANCE

The chi-square value is often used to judge the significance of population variance i.e., we can use the test to judge if a random sample has been drawn from a normal population with mean (μ) and with a specified variance (σ_p^2). The test is based on χ^2 -distribution. Such a distribution we encounter when we deal with collections of values that involve adding up squares. Variances of samples require us to add a collection of squared quantities and, thus, have distributions that are related to χ^2 -distribution. If we take each one of a collection of sample variances, divided them by the known population variance and multiply these quotients by $(n - 1)$, where n means the number of items in

the sample, we shall obtain a χ^2 -distribution. Thus, $\frac{\sigma_s^2}{\sigma_p^2}(n - 1) = \frac{\sigma_s^2}{\sigma_p^2}$ (d.f.) would have the same distribution as χ^2 -distribution with $(n - 1)$ degrees of freedom.

* See Chapter 12 Testing of Hypotheses-II for more details.

¹ Neil R. Ullman, *Elementary Statistics—An Applied Approach*, p. 234.

The χ^2 -distribution is not symmetrical and all the values are positive. For making use of this distribution, one is required to know the degrees of freedom since for different degrees of freedom we have different curves. The smaller the number of degrees of freedom, the more skewed is the distribution which is illustrated in Fig. 10.1:

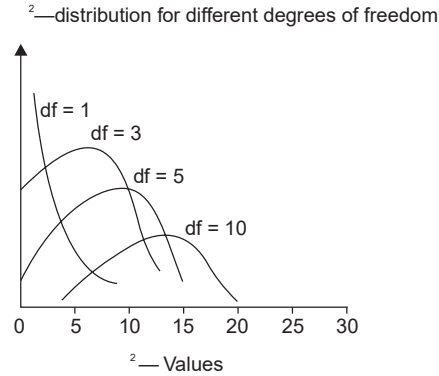


Fig. 10.1

Table given in the Appendix gives selected critical values of χ^2 for the different degrees of freedom. χ^2 -values are the quantities indicated on the x-axis of the above diagram and in the table are areas below that value.

In brief, when we have to use chi-square as a test of population variance, we have to work out the value of χ^2 to test the null hypothesis (viz., $H_0: \sigma_s^2 = \sigma_p^2$) as under:

$$\chi^2 = \frac{\sigma_s^2}{\sigma_p^2}(n - 1)$$

where σ_s^2 = variance of the sample;

σ_p^2 = variance of the population;

$(n - 1)$ = degrees of freedom, n being the number of items in the sample.

Then by comparing the calculated value with the table value of χ^2 for $(n - 1)$ degrees of freedom at a given level of significance, we may either accept or reject the null hypothesis. If the calculated value of χ^2 is less than the table value, the null hypothesis is accepted, but if the calculated value is equal or greater than the table value, the hypothesis is rejected. All this can be made clear by an example.

Illustration 1

Weight of 10 students is as follows:

S. No.	1	2	3	4	5	6	7	8	9	10
Weight (kg.)	38	40	45	53	47	43	55	48	52	49

Can we say that the variance of the distribution of weight of all students from which the above sample of 10 students was drawn is equal to 20 kgs? Test this at 5 per cent and 1 per cent level of significance.

Solution: First of all we should work out the variance of the sample data or σ_s^2 and the same has been worked out as under:

Table 10.1

S. No.	X_i (Weight in kgs.)	$(X_i - \bar{X})$	$(X_i - \bar{X})^2$
1	38	-9	81
2	40	-7	49
3	45	-2	04
4	53	+6	36
5	47	+0	00
6	43	-4	16
7	55	+8	64
8	48	+1	01
9	52	+5	25
10	49	+2	04
$n = 10$	$\sum X_i = 470$	$\sum (X_i - \bar{X})^2 = 280$	

$$\bar{X} = \frac{\sum X_i}{n} = \frac{470}{10} = 47 \text{ kgs.}$$

$$\therefore \sigma_s = \sqrt{\frac{\sum (X_i - \bar{X})^2}{n - 1}} = \sqrt{\frac{280}{10 - 1}} = \sqrt{31.11}$$

or $\sigma_s^2 = 31.11.$

Let the null hypothesis be $H_0 : \sigma_p^2 = \sigma_s^2$. In order to test this hypothesis we work out the χ^2 value as under:

$$\chi^2 = \frac{\sigma_s^2}{\sigma_p^2} (n - 1)$$

$$= \frac{31.11}{20}(10 - 1) = 13.999.$$

Degrees of freedom in the given case is $(n - 1) = (10 - 1) = 9$. At 5 per cent level of significance the table value of $\chi^2 = 16.92$ and at 1 per cent level of significance, it is 21.67 for 9 d.f. and both these values are greater than the calculated value of χ^2 which is 13.999. Hence we accept the null hypothesis and conclude that the variance of the given distribution can be taken as 20 kgs at 5 per cent as also at 1 per cent level of significance. In other words, the sample can be said to have been taken from a population with variance 20 kgs.

Illustration 2

A sample of 10 is drawn randomly from a certain population. The sum of the squared deviations from the mean of the given sample is 50. Test the hypothesis that the variance of the population is 5 at 5 per cent level of significance.

Solution: Given information is

$$n = 10$$

$$\Sigma(X_i - \bar{X})^2 = 50$$

$$\therefore \sigma_s^2 = \frac{\Sigma(X_i - \bar{X})^2}{n - 1} = \frac{50}{9}$$

Take the null hypothesis as $H_0: \sigma_p^2 = \sigma_s^2$. In order to test this hypothesis, we work out the χ^2 value as under:

$$\chi^2 = \frac{\sigma_s^2}{\sigma_p^2}(n - 1) = \frac{\frac{50}{9}}{5}(10 - 1) = \frac{50}{9} \times \frac{1}{5} \times \frac{9}{1} = 10$$

Degrees of freedom $= (10 - 1) = 9$.

The table value of χ^2 at 5 per cent level for 9 d.f. is 16.92. The calculated value of χ^2 is less than this table value, so we accept the null hypothesis and conclude that the variance of the population is 5 as given in the question.

CHI-SQUARE AS A NON-PARAMETRIC TEST

Chi-square is an important non-parametric test and as such no rigid assumptions are necessary in respect of the type of population. We require only the degrees of freedom (implicitly of course the size of the sample) for using this test. As a non-parametric test, chi-square can be used (i) as a test of goodness of fit and (ii) as a test of independence.

As a test of goodness of fit, χ^2 test enables us to see how well does the assumed theoretical distribution (such as Binomial distribution, Poisson distribution or Normal distribution) fit to the observed data. When some theoretical distribution is fitted to the given data, we are always interested in knowing as to how well this distribution fits with the observed data. The chi-square test can give answer to this. If the calculated value of χ^2 is less than the table value at a certain level of significance, the fit is considered to be a good one which means that the divergence between the observed and expected frequencies is attributable to fluctuations of sampling. But if the calculated value of χ^2 is greater than its table value, the fit is not considered to be a good one.

As a test of independence, χ^2 test enables us to explain whether or not two attributes are associated. For instance, we may be interested in knowing whether a new medicine is effective in controlling fever or not, χ^2 test will helps us in deciding this issue. In such a situation, we proceed with the null hypothesis that the two attributes (viz., new medicine and control of fever) are independent which means that new medicine is not effective in controlling fever. On this basis we first calculate the expected frequencies and then work out the value of χ^2 . If the calculated value of χ^2 is less than the table value at a certain level of significance for given degrees of freedom, we conclude that null hypothesis stands which means that the two attributes are independent or not associated (i.e., the new medicine is not effective in controlling the fever). But if the calculated value of χ^2 is greater than its table value, our inference then would be that null hypothesis does not hold good which means the two attributes are associated and the association is not because of some chance factor but it exists in reality (i.e., the new medicine is effective in controlling the fever and as such may be prescribed). It may, however, be stated here that χ^2 is not a measure of the degree of relationship or the form of relationship between two attributes, but is simply a technique of judging the significance of such association or relationship between two attributes.

In order that we may apply the chi-square test either as a test of goodness of fit or as a test to judge the significance of association between attributes, it is necessary that the observed as well as theoretical or expected frequencies must be grouped in the same way and the theoretical distribution must be adjusted to give the same total frequency as we find in case of observed distribution. χ^2 is then calculated as follows:

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

where

O_{ij} = observed frequency of the cell in i th row and j th column.

E_{ij} = expected frequency of the cell in i th row and j th column.

If two distributions (observed and theoretical) are exactly alike, $\chi^2 = 0$; but generally due to sampling errors, χ^2 is not equal to zero and as such we must know the sampling distribution of χ^2 so that we may find the probability of an observed χ^2 being given by a random sample from the hypothetical universe. Instead of working out the probabilities, we can use ready table which gives probabilities for given values of χ^2 . Whether or not a calculated value of χ^2 is significant can be

ascertained by looking at the tabulated values of χ^2 for given degrees of freedom at a certain level of significance. If the calculated value of χ^2 is equal to or exceeds the table value, the difference between the observed and expected frequencies is taken as significant, but if the table value is more than the calculated value of χ^2 , then the difference is considered as insignificant i.e., considered to have arisen as a result of chance and as such can be ignored.

As already stated, degrees of freedom* play an important part in using the chi-square distribution and the test based on it, one must correctly determine the degrees of freedom. If there are 10 frequency classes and there is one independent constraint, then there are $(10 - 1) = 9$ degrees of freedom. Thus, if 'n' is the number of groups and one constraint is placed by making the totals of observed and expected frequencies equal, the d.f. would be equal to $(n - 1)$. In the case of a contingency table (i.e., a table with 2 columns and 2 rows or a table with two columns and more than two rows or a table with two rows but more than two columns or a table with more than two rows and more than two columns), the d.f. is worked out as follows:

$$\text{d.f.} = (c - 1)(r - 1)$$

where 'c' means the number of columns and 'r' means the number of rows.

CONDITIONS FOR THE APPLICATION OF χ^2 TEST

The following conditions should be satisfied before χ^2 test can be applied:

- (i) Observations recorded and used are collected on a random basis.
- (ii) All the items in the sample must be independent.
- (iii) No group should contain very few items, say less than 10. In case where the frequencies are less than 10, regrouping is done by combining the frequencies of adjoining groups so that the new frequencies become greater than 10. Some statisticians take this number as 5, but 10 is regarded as better by most of the statisticians.
- (iv) The overall number of items must also be reasonably large. It should normally be at least 50, howsoever small the number of groups may be.
- (v) The constraints must be linear. Constraints which involve linear equations in the cell frequencies of a contingency table (i.e., equations containing no squares or higher powers of the frequencies) are known as linear constraints.

STEPS INVOLVED IN APPLYING CHI-SQUARE TEST

The various steps involved are as follows:

* For d.f. greater than 30, the distribution of $\sqrt{2\chi^2}$ approximates the normal distribution wherein the mean of $\sqrt{2\chi^2}$ distribution is $\sqrt{2\text{d.f.} - 1}$ and the standard deviation = 1. Accordingly, when d.f. exceeds 30, the quantity $\left[\sqrt{2\chi^2} - \sqrt{2\text{d.f.} - 1} \right]$ may be used as a normal variate with unit variance, i.e.,

$$z_\alpha = \sqrt{2\chi^2} - \sqrt{2\text{d.f.} - 1}$$

- (i) First of all calculate the expected frequencies on the basis of given hypothesis or on the basis of null hypothesis. Usually in case of a 2×2 or any contingency table, the expected frequency for any given cell is worked out as under:

$$\text{Expected frequency of any cell} = \left[\frac{(\text{Row total for the row of that cell}) \times (\text{Column total for the column of that cell})}{(\text{Grand total})} \right]$$

- (ii) Obtain the difference between observed and expected frequencies and find out the squares of such differences i.e., calculate $(O_{ij} - E_{ij})^2$.
- (iii) Divide the quantity $(O_{ij} - E_{ij})^2$ obtained as stated above by the corresponding expected frequency to get $(O_{ij} - E_{ij})^2/E_{ij}$ and this should be done for all the cell frequencies or the group frequencies.

- (iv) Find the summation of $(O_{ij} - E_{ij})^2/E_{ij}$ values or what we call $\sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$. This is the required χ^2 value.

The χ^2 value obtained as such should be compared with relevant table value of χ^2 and then inference be drawn as stated above.

We now give few examples to illustrate the use of χ^2 test.

Illustration 3

A die is thrown 132 times with following results:

Number turned up	1	2	3	4	5	6
Frequency	16	20	25	14	29	28

Is the die unbiased?

Solution: Let us take the hypothesis that the die is unbiased. If that is so, the probability of obtaining any one of the six numbers is $1/6$ and as such the expected frequency of any one number coming upward is $132 \times 1/6 = 22$. Now we can write the observed frequencies along with expected frequencies and work out the value of χ^2 as follows:

Table 10.2

No. turned up	Observed frequency O_i	Expected frequency E_i	$(O_i - E_i)$	$(O_i - E_i)^2$	$(O_i - E_i)^2/E_i$
1	16	22	-6	36	36/22
2	20	22	-2	4	4/22
3	25	22	3	9	9/22
4	14	22	-8	64	64/22
5	29	22	7	49	49/22
6	28	22	6	36	36/22

$$\therefore \sum [(O_i - E_i)^2 / E_i] = 9.$$

Hence, the calculated value of $\chi^2 = 9$.

\therefore Degrees of freedom in the given problem is

$$(n - 1) = (6 - 1) = 5.$$

The table value* of χ^2 for 5 degrees of freedom at 5 per cent level of significance is 11.071. Comparing calculated and table values of χ^2 , we find that calculated value is less than the table value and as such could have arisen due to fluctuations of sampling. The result, thus, supports the hypothesis and it can be concluded that the die is unbiased.

Illustration 4

Find the value of χ^2 for the following information:

Class	A	B	C	D	E
Observed frequency	8	29	44	15	4
Theoretical (or expected) frequency	7	24	38	24	7

Solution: Since some of the frequencies less than 10, we shall first re-group the given data as follows and then will work out the value of χ^2 :

Table 10.3

Class	Observed frequency O_i	Expected frequency E_i	$O_i - E_i$	$(O_i - E_i)^2 / E_i$
A and B	$(8 + 29) = 37$	$(7 + 24) = 31$	6	36/31
C	44	38	6	36/38
D and E	$(15 + 4) = 19$	$(24 + 7) = 31$	-12	144/31

$$\therefore \chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} = 6.76 \text{ app.}$$

Illustration 5

Genetic theory states that children having one parent of blood type A and the other of blood type B will always be of one of three types, A, AB, B and that the proportion of three types will on an average be as 1 : 2 : 1. A report states that out of 300 children having one A parent and B parent, 30 per cent were found to be types A, 45 per cent per cent type AB and remainder type B. Test the hypothesis by χ^2 test.

Solution: The observed frequencies of type A, AB and B is given in the question are 90, 135 and 75 respectively.

*Table No. 3 showing some critical values of χ^2 for specified degrees of freedom has been given in Appendix at the end of the book.

The expected frequencies of type A , AB and B (as per the genetic theory) should have been 75, 150 and 75 respectively.

We now calculate the value of χ^2 as follows:

Table 10.4

Type	Observed frequency O_i	Expected frequency E_i	$(O_i - E_i)$	$(O_i - E_i)^2$	$(O_i - E_i)^2/E_i$
A	90	75	15	225	$225/75 = 3$
AB	135	150	-15	225	$225/150 = 1.5$
B	75	75	0	0	$0/75 = 0$

$$\therefore \chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} = 3 + 1.5 + 0 = 4.5$$

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

Table value of χ^2 for 2 d.f. at 5 per cent level of significance is 5.991.

The calculated value of χ^2 is 4.5 which is less than the table value and hence can be ascribed to have taken place because of chance. This supports the theoretical hypothesis of the genetic theory that on an average type A , AB and B stand in the proportion of 1 : 2 : 1.

Illustration 6

The table given below shows the data obtained during outbreak of smallpox:

	Attacked	Not attacked	Total
Vaccinated	31	469	500
Not vaccinated	185	1315	1500
Total	216	1784	2000

Test the effectiveness of vaccination in preventing the attack from smallpox. Test your result with the help of χ^2 at 5 per cent level of significance.

Solution: Let us take the hypothesis that vaccination is not effective in preventing the attack from smallpox i.e., vaccination and attack are independent. On the basis of this hypothesis, the expected frequency corresponding to the number of persons vaccinated and attacked would be:

$$\text{Expectation of } (AB) = \frac{(A) \times (B)}{N}$$

when A represents vaccination and B represents attack.

$$\therefore \begin{aligned} (A) &= 500 \\ (B) &= 216 \\ N &= 2000 \end{aligned}$$

$$\text{Expectation of } (AB) = \frac{500 \times 216}{2000} = 54$$

Now using the expectation of (AB) , we can write the table of expected values as follows:

	Attacked: B	Not attacked: b	Total
Vaccinated: A	$(AB) = 54$	$(Ab) = 446$	500
Not vaccinated: a	$(aB) = 162$	$(ab) = 1338$	1500
Total	216	1784	2000

Table 10.5: Calculation of Chi-Square

Group	Observed frequency O_{ij}	Expected frequency E_{ij}	$(O_{ij} - E_{ij})$	$(O_{ij} - E_{ij})^2$	$(O_{ij} - E_{ij})^2/E_{ij}$
AB	31	54	-23	529	$529/54 = 9.796$
Ab	469	446	+23	529	$529/44 = 1.186$
aB	158	162	+23	529	$529/162 = 3.265$
ab	1315	1338	-23	529	$529/1338 = 0.395$

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = 14.642$$

\therefore Degrees of freedom in this case $= (r - 1)(c - 1) = (2 - 1)(2 - 1) = 1$.

The table value of χ^2 for 1 degree of freedom at 5 per cent level of significance is 3.841. The calculated value of χ^2 is much higher than this table value and hence the result of the experiment does not support the hypothesis. We can, thus, conclude that vaccination is effective in preventing the attack from smallpox.

Illustration 7

Two research workers classified some people in income groups on the basis of sampling studies. Their results are as follows:

Investigators	Income groups			Total
	Poor	Middle	Rich	
A	160	30	10	200
B	140	120	40	300
Total	300	150	50	500

Show that the sampling technique of at least one research worker is defective.

Solution: Let us take the hypothesis that the sampling techniques adopted by research workers are similar (i.e., there is no difference between the techniques adopted by research workers). This being so, the expectation of *A* investigator classifying the people in

$$(i) \text{ Poor income group} = \frac{200 \times 300}{500} = 120$$

$$(ii) \text{ Middle income group} = \frac{200 \times 150}{500} = 60$$

$$(iii) \text{ Rich income group} = \frac{200 \times 50}{500} = 20$$

Similarly the expectation of *B* investigator classifying the people in

$$(i) \text{ Poor income group} = \frac{300 \times 300}{500} = 180$$

$$(ii) \text{ Middle income group} = \frac{300 \times 150}{500} = 90$$

$$(iii) \text{ Rich income group} = \frac{300 \times 50}{500} = 30$$

We can now calculate value of χ^2 as follows:

Table 10.6

Groups	Observed frequency O_{ij}	Expected frequency E_{ij}	$O_{ij} - E_{ij}$	$(O_{ij} - E_{ij})^2 E_{ij}$
<i>Investigator A</i>				
classifies people as poor	160	120	40	1600/120 = 13.33
classifies people as middle class people	30	60	-30	900/60 = 15.00
classifies people as rich	10	20	-10	100/20 = 5.00
<i>Investigator B</i>				
classifies people as poor	140	180	-40	1600/180 = 8.88
classifies people as middle class people	120	90	30	900/90 = 10.00
classifies people as rich	40	30	10	100/30 = 3.33

Hence,
$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = 55.54$$

\therefore Degrees of freedom $= (c - 1)(r - 1)$
 $= (3 - 1)(2 - 1) = 2.$

The table value of χ^2 for two degrees of freedom at 5 per cent level of significance is 5.991.

The calculated value of χ^2 is much higher than this table value which means that the calculated value cannot be said to have arisen just because of chance. It is significant. Hence, the hypothesis does not hold good. This means that the sampling techniques adopted by two investigators differ and are not similar. Naturally, then the technique of one must be superior than that of the other.

Illustration 8

Eight coins were tossed 256 times and the following results were obtained:

<i>Numbers of heads</i>	0	1	2	3	4	5	6	7	8
<i>Frequency</i>	2	6	30	52	67	56	32	10	1

Are the coins biased? Use χ^2 test.

Solution: Let us take the hypothesis that the coins are not biased. If that is so, the probability of any one coin falling with head upward is $1/2$ and with tail upward is $1/2$ and it remains the same whatever be the number of throws. In such a case the expected values of getting 0, 1, 2, ... heads in a single throw in 256 throws of eight coins will be worked out as follows*.

Table 10.7

<i>Events or No. of heads</i>	<i>Expected frequencies</i>
0	${}^8C_0 \left(\frac{1}{2}\right)^0 \left(\frac{1}{2}\right)^8 \times 256 = 1$
1	${}^8C_1 \left(\frac{1}{2}\right)^1 \left(\frac{1}{2}\right)^7 \times 256 = 8$
2	${}^8C_2 \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^6 \times 256 = 28$

contd.

* The probabilities of random variable i.e., various possible events have been worked out on the binomial principle viz., through the expansion of $(p + q)^n$ where $p = 1/2$ and $q = 1/2$ and $n = 8$ in the given case. The expansion of the term ${}^nC_r p^r q^{n-r}$ has given the required probabilities which have been multiplied by 256 to obtain the expected frequencies.

<i>Events or No. of heads</i>	<i>Expected frequencies</i>
3	${}^8C_3 \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^5 \times 256 = 56$
4	${}^8C_4 \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^4 \times 256 = 70$
5	${}^8C_5 \left(\frac{1}{2}\right)^5 \left(\frac{1}{2}\right)^3 \times 256 = 56$
6	${}^8C_6 \left(\frac{1}{2}\right)^6 \left(\frac{1}{2}\right)^2 \times 256 = 28$
7	${}^8C_7 \left(\frac{1}{2}\right)^7 \left(\frac{1}{2}\right)^1 \times 256 = 8$
8	${}^8C_8 \left(\frac{1}{2}\right)^8 \left(\frac{1}{2}\right)^0 \times 256 = 1$

The value of χ^2 can be worked out as follows:

Table 10.8

<i>No. of heads</i>	<i>Observed frequency O_i</i>	<i>Expected frequency E_i</i>	$O_i - E_i$	$(O_i - E_i)^2/E_i$
0	2	1	1	1/1 = 1.00
1	6	8	-2	4/8 = 0.50
2	30	28	2	4/28 = 0.14
3	52	56	-4	16/56 = 0.29
4	67	70	-3	9/70 = 0.13
5	56	56	0	0/56 = 0.00
6	32	28	4	16/28 = 0.57
7	10	8	2	4/8 = 0.50
8	1	1	0	0/1 = 0.00

$$\therefore \chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} = 3.13$$

\therefore Degrees of freedom = $(n - 1) = (9 - 1) = 8$

The table value of χ^2 for eight degrees of freedom at 5 per cent level of significance is 15.507.

The calculated value of χ^2 is much less than this table and hence it is insignificant and can be ascribed due to fluctuations of sampling. The result, thus, supports the hypothesis and we may say that the coins are not biased.

ALTERNATIVE FORMULA

There is an alternative method of calculating the value of χ^2 in the case of a (2×2) table. If we write the cell frequencies and marginal totals in case of a (2×2) table thus,

a	b	$(a + b)$
c	d	$(c + d)$
$(a + c)$	$(b + d)$	N

then the formula for calculating the value of χ^2 will be stated as follows:

$$\chi^2 = \frac{(ad - bc)^2 \cdot N}{(a + c)(b + d)(a + b)(c + d)}$$

where N means the total frequency, ad means the larger cross product, bc means the smaller cross product and $(a + c)$, $(b + d)$, $(a + b)$, and $(c + d)$ are the marginal totals. The alternative formula is rarely used in finding out the value of chi-square as it is not applicable uniformly in all cases but can be used only in a (2×2) contingency table.

YATES' CORRECTION

F. Yates has suggested a correction for continuity in χ^2 value calculated in connection with a (2×2) table, particularly when cell frequencies are small (since no cell frequency should be less than 5 in any case, though 10 is better as stated earlier) and χ^2 is just on the significance level. The correction suggested by Yates is popularly known as Yates' correction. It involves the reduction of the deviation of observed from expected frequencies which of course reduces the value of χ^2 . The rule for correction is to adjust the observed frequency in each cell of a (2×2) table in such a way as to reduce the deviation of the observed from the expected frequency for that cell by 0.5, but this adjustment is made in all the cells without disturbing the marginal totals. The formula for finding the value of χ^2 after applying Yates' correction can be stated thus:

$$\chi^2(\text{corrected}) = \frac{N \cdot (|ad - bc| - 0.5N)^2}{(a + b)(c + d)(a + c)(b + d)}$$

In case we use the usual formula for calculating the value of chi-square viz.,

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}},$$

then Yates' correction can be applied as under:

$$\chi^2(\text{corrected}) = \frac{[|O_1 - E_1| - 0.5]^2}{E_1} + \frac{[|O_2 - E_2| - 0.5]^2}{E_2} + \dots$$

It may again be emphasised that Yates' correction is made only in case of (2×2) table and that too when cell frequencies are small.

Illustration 9

The following information is obtained concerning an investigation of 50 ordinary shops of small size:

	Shops		Total
	In towns	In villages	
Run by men	17	18	35
Run by women	3	12	15
Total	20	30	50

Can it be inferred that shops run by women are relatively more in villages than in towns? Use χ^2 test.

Solution: Take the hypothesis that there is no difference so far as shops run by men and women in towns and villages. With this hypothesis the expectation of shops run by men in towns would be:

$$\text{Expectation of } (AB) = \frac{(A) \times (B)}{N}$$

where A = shops run by men

B = shops in towns

$(A) = 35; (B) = 20$ and $N = 50$

$$\text{Thus, expectation of } (AB) = \frac{35 \times 20}{50} = 14$$

Hence, table of expected frequencies would be

	<i>Shops in towns</i>	<i>Shops in villages</i>	<i>Total</i>
Run by men	14 (<i>AB</i>)	21 (<i>Ab</i>)	35
Run by women	6 (<i>aB</i>)	9 (<i>ab</i>)	15
Total	20	30	50

Calculation of χ^2 value:

Table 10.9

<i>Groups</i>	<i>Observed frequency</i> O_{ij}	<i>Expected frequency</i> E_{ij}	$(O_{ij} - E_{ij})$	$(O_{ij} - E_{ij})^2/E_{ij}$
(<i>AB</i>)	17	14	3	9/14=0.64
(<i>Ab</i>)	18	21	-3	9/21=0.43
(<i>aB</i>)	3	6	-3	9/6=1.50
(<i>ab</i>)	12	9	3	9/9=1.00

$$\therefore \chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = 3.57$$

As one cell frequency is only 3 in the given 2×2 table, we also work out χ^2 value applying Yates' correction and this is as under:

$$\begin{aligned} \chi^2 (\text{corrected}) &= \frac{[|17 - 14| - 0.5]^2}{14} + \frac{[|18 - 21| - 0.5]^2}{21} + \frac{[|3 - 6| - 0.5]^2}{6} + \frac{[|12 - 9| - 0.5]^2}{9} \\ &= \frac{(2.5)^2}{14} + \frac{(2.5)^2}{21} + \frac{(2.5)^2}{6} + \frac{(2.5)^2}{9} \\ &= 0.446 + 0.298 + 1.040 + 0.694 \\ &= 2.478 \end{aligned}$$

$$\therefore \text{Degrees of freedom} = (c - 1)(r - 1) = (2 - 1)(2 - 1) = 1$$

Table value of χ^2 for one degree of freedom at 5 per cent level of significance is 3.841. The calculated value of χ^2 by both methods (i.e., before correction and after Yates' correction) is less than its table value. Hence the hypothesis stands. We can conclude that there is no difference between shops run by men and women in villages and towns.

Additive property: An important property of χ^2 is its additive nature. This means that several values of χ^2 can be added together and if the degrees of freedom are also added, this number gives the degrees of freedom of the total value of χ^2 . Thus, if a number of χ^2 values have been obtained

from a number of samples of similar data, then because of the additive nature of χ^2 we can combine the various values of χ^2 by just simply adding them. Such addition of various values of χ^2 gives one value of χ^2 which helps in forming a better idea about the significance of the problem under consideration. The following example illustrates the additive property of χ^2 .

Illustration 10

The following values of χ^2 from different investigations carried to examine the effectiveness of a recently invented medicine for checking malaria are obtained:

Investigation	χ^2	d.f.
1	2.5	1
2	3.2	1
3	4.1	1
4	3.7	1
5	4.5	1

What conclusion would you draw about the effectiveness of the new medicine on the basis of the five investigations taken together?

Solution: By adding all the values of χ^2 , we obtain a value equal to 18.0. Also by adding the various d.f., as given in the question, we obtain the value 5. We can now state that the value of χ^2 for 5 degrees of freedom (when all the five investigations are taken together) is 18.0.

Let us take the hypothesis that the new medicine is not effective. The table value of χ^2 for 5 degrees of freedom at 5 per cent level of significance is 11.070. But our calculated value is higher than this table value which means that the difference is significant and is not due to chance. As such the hypothesis is rejected and it can be concluded that the new medicine is effective in checking malaria.

CONVERSION OF CHI-SQUARE INTO PHI COEFFICIENT (ϕ)

Since χ^2 does not by itself provide an estimate of the magnitude of association between two attributes, any obtained χ^2 value may be converted into Phi coefficient (symbolized as ϕ) for the purpose. In other words, chi-square tells us about the significance of a relation between variables; it provides no answer regarding the magnitude of the relation. This can be achieved by computing the Phi coefficient, which is a non-parametric measure of coefficient of correlation, as under:

$$\phi = \sqrt{\frac{\chi^2}{N}}$$

CONVERSION OF CHI-SQUARE INTO COEFFICIENT OF CONTINGENCY (C)

Chi-square value may also be converted into coefficient of contingency, especially in case of a contingency table of higher order than 2×2 table to study the magnitude of the relation or the degree of association between two attributes, as shown below:

$$C = \sqrt{\frac{\chi^2}{\chi^2 + N}}$$

While finding out the value of C we proceed on the assumption of null hypothesis that the two attributes are independent and exhibit no association. Coefficient of contingency is also known as coefficient of Mean Square contingency. This measure also comes under the category of non-parametric measure of relationship.

IMPORTANT CHARACTERISTICS OF χ^2 TEST

- (i) This test (as a non-parametric test) is based on frequencies and not on the parameters like mean and standard deviation.
- (ii) The test is used for testing the hypothesis and is not useful for estimation.
- (iii) This test possesses the additive property as has already been explained.
- (iv) This test can also be applied to a complex contingency table with several classes and as such is a very useful test in research work.
- (v) This test is an important non-parametric test as no rigid assumptions are necessary in regard to the type of population, no need of parameter values and relatively less mathematical details are involved.

CAUTION IN USING χ^2 TEST

The chi-square test is no doubt a most frequently used test, but its correct application is equally an uphill task. It should be borne in mind that the test is to be applied only when the individual observations of sample are independent which means that the occurrence of one individual observation (event) has no effect upon the occurrence of any other observation (event) in the sample under consideration. Small theoretical frequencies, if these occur in certain groups, should be dealt with under special care. The other possible reasons concerning the improper application or misuse of this test can be (i) neglect of frequencies of non-occurrence; (ii) failure to equalise the sum of observed and the sum of the expected frequencies; (iii) wrong determination of the degrees of freedom; (iv) wrong computations, and the like. The researcher while applying this test must remain careful about all these things and must thoroughly understand the rationale of this important test before using it and drawing inferences in respect of his hypothesis.

Questions

1. What is Chi-square test? Explain its significance in statistical analysis.
2. Write short notes on the following:
 - (i) Additive property of Chi-square;
 - (ii) Chi-square as a test of 'goodness of fit';
 - (iii) Precautions in applying Chi-square test;
 - (iv) Conditions for applying Chi-square test.
3. An experiment was conducted to test the efficacy of chloromycetin in checking typhoid. In a certain hospital chloromycetin was given to 285 out of the 392 patients suffering from typhoid. The number of typhoid cases were as follows:

	<i>Typhoid</i>	<i>No Typhoid</i>	<i>Total</i>
Chloromycetin	35	250	285
No chloromycetin	50	57	107
Total	85	307	392

With the help of χ^2 , test the effectiveness of chloromycetin in checking typhoid.

(The χ^2 value at 5 per cent level of significance for one degree of freedom is 3.841).

(*M. Com., Rajasthan University, 1966*)

4. On the basis of information given below about the treatment of 200 patients suffering from a disease, state whether the new treatment is comparatively superior to the conventional treatment.

<i>Treatment</i>	<i>No. of patients</i>	
	<i>Favourable Response</i>	<i>No Response</i>
New	60	20
Conventional	70	50

For drawing your inference, use the value of χ^2 for one degree of freedom at the 5 per cent level of significance, viz., 3.84.

5. 200 digits were chosen at random from a set of tables. The frequencies of the digits were:

Digit	0	1	2	3	4	5	6	7	8	9
Frequency	18	19	23	21	16	25	22	20	21	15

Calculate χ^2 .

6. Five dice were thrown 96 times and the number of times 4, 5, or 6 was thrown were

Number of dice throwing

4, 5 or 6	5	4	3	2	1	0
Frequency	8	18	35	24	10	1

Find the value of Chi-square.

7. Find Chi-square from the following information:

Condition of child	Condition of home		Total
	Clean	Dirty	
Clean	70	50	120
Fairly clean	80	20	100
Dirty	35	45	80
Total	185	115	300

State whether the two attributes viz., condition of home and condition of child are independent (Use Chi-square test for the purpose).

8. In a certain cross the types represented by XY , Xy , xY and xy are expected to occur in a 9 : 5 : 4 : 2 ratio. The actual frequencies were:

XY	Xy	xY	xy
180	110	60	50

Test the goodness of fit of observation to theory.

9. The normal rate of infection for a certain disease in cattle is known to be 50 per cent. In an experiment with seven animals injected with a new vaccine it was found that none of the animals caught infection. Can the evidence be regarded as conclusive (at 1 per cent level of significance) to prove the value of the new vaccine?
10. Result of throwing die were recorded as follows:
- | | | | | | | |
|------------------------|----|----|----|----|----|----|
| Number falling upwards | 1 | 2 | 3 | 4 | 5 | 6 |
| Frequency | 27 | 33 | 31 | 29 | 30 | 24 |
- Is the die unbiased? Answer on the basis of Chi-square test.
11. The Theory predicts the proportion of beans, in the four groups A, B, C and D should be 9 : 3 : 3 : 1. In an experiment among 1600 beans, the number in the four groups were 882, 313, 287 and 118. Does the experimental result support the theory? Apply χ^2 test.

(M.B.A., Delhi University, 1975)

12. You are given a sample of 150 observations classified by two attributes A and B as follows:

	A_1	A_2	A_3	Total
B_1	40	25	15	80
B_2	11	26	8	45
B_3	9	9	7	25
Total	60	60	30	150

Use the χ^2 test to examine whether A and B are associated.

(M.A. Eco., Patiala University, 1975)

13. A survey of 320 families with five children each revealed the following distribution:

No. of boys	5	4	3	2	1	0
No. of girls	0	1	2	3	4	5
No. of families	14	56	110	88	40	12

Is this distribution consistent with the hypothesis that male and female births are equally probable? Apply Chi-square test.

14. What is Yates' correction? Find the value of Chi-square applying Yates' correction to the following data:

	<i>Passed</i>	<i>Failed</i>	<i>Total</i>
Day classes	10	20	30
Evening classes	4	66	70
Total	14	86	100

Also state whether the association, if any, between passing in the examination and studying in day classes is significant using Chi-square test.

15. (a) 1000 babies were born during a certain week in a city of which 600 were boys and 400 girls. Use χ^2 test to examine the correctness of the hypothesis that the sex-ratio is 1 : 1 in newly born babies.
 (b) The percentage of smokers in a certain city was 90. A random sample of 100 persons was selected in which 85 persons were found to be smokers. Is the sample proportion significantly different from the proportion of smokers in the city? Answer on the basis of Chi-square test.
16. A college is running post-graduate classes in five subjects with equal number of students. The total number of absentees in these five classes is 75. Test the hypothesis that these classes are alike in absenteeism if the actual absentees in each are as follows:

History	= 19
Philosophy	= 18
Economics	= 15
Commerce	= 12
Chemistry	= 11

(M.Phil. (EAFM) Exam. Raj. Uni., 1978)

17. The number of automobile accidents per week in a certain community were as follows:

12, 8, 20, 2, 14, 10, 15, 6, 9, 4

Are these frequencies in agreement with the belief that accident conditions were the same during the 10 week period under consideration?

18. A certain chemical plant processes sea water to collect sodium chloride and magnesium. From scientific analysis, sea water is known to contain sodium chloride, magnesium and other elements in the ratio of 62 : 4 : 34. A sample of 200 tons of sea water has resulted in 130 tons of sodium chloride and 6 tons of magnesium. Are these data consistent with the scientific model at 5 per cent level of significance?
19. An oil company has explored three different areas for possible oil reserves. The results of the test were as given below:

	Area			Total
	A	B	C	
Strikes	7	10	8	25
Dry holes	10	18	9	37
Total number of test wells	17	28	17	62

Do the three areas have the same potential, at the 10 per cent level of significance?

20. While conducting an air traffic study, a record was made of the number of aircraft arrivals, at a certain airport, during 250 half hour time intervals. The following tables gives the observed number of periods in which there were 0, 1, 2, 3, 4, or more arrivals as well as the expected number of such periods if arrivals per half hour have a Poisson distribution $\lambda = 2$. Does this Poisson distribution describe the observed arrivals at 5 per cent level of significance.

Number of observed arrivals (per half hour)	Number of periods observed	Number of periods expected (Poisson, $\lambda = 2$)
0	47	34
1	56	68
2	71	68
3	44	45
4 or more	32	35

21. A marketing researcher interested in the business publication reading habits of purchasing agents has assembled the following data:

Business Publication Preferences (First Choice Mentions)

Business Publication	Frequency of first choice
A	35
B	30
C	45
D	55

- (i) Test the null hypothesis ($\alpha = 0.05$) that there are no differences among frequencies of first choice of tested publications.
- (ii) If the choice of A and C and that of B and D are aggregated, test the null hypothesis at $\alpha = 0.05$ that there are no differences.
22. A group of 150 College students were asked to indicate their most liked film star from among six different well known film actors viz., A, B, C, D, E and F in order to ascertain their relative popularity. The observed frequency data were as follows:

Actors	A	B	C	D	E	F	Total
Frequencies	24	20	32	25	28	21	150

Test at 5 per cent whether all actors are equally popular.

23. For the data in question 12, find the coefficient of contingency to measure the magnitude of relationship between A and B .
24. (a) What purpose is served by calculating the Phi coefficient (ϕ)? Explain.
- (b) If $\chi^2 = 16$ and $N = 4$, find the value of Phi coefficient.