

Probability and Statistics
Prof. Dr. Somesh Kumar
Department of Mathematics
Indian Institute of Technology, Kharagpur

Module No. #01
Lecture No. #39
Testing of Hypothesis-VII

In the situations so far, for testing of hypothesis problems, we have considered that the random sample comes from a certain population. So, we have assumed the form of population to be say, normal or say exponential. And then, we want to test about the parameters of that population. For example, we want to test about the mean of a normal distribution, we want to test about the variance of a normal distribution or we want to test about the scale parameter of an exponential distribution. So, here the situation is that, we assume that the form of the distribution is known to us. Only the parameters of the population are not known. However, there are other situations where we have a data and we want to know that, from which type of population that data has come. And so, that means, we may be liking to estimate the population, that means, the distribution or we want to test about the distribution. So, here we will talk about the testing, that a particular distribution is there, say capital F is equal to sum F naught. So, for this situation an approximate test has been proposed which is called chi square test for goodness of fit.

(Refer Slide Time: 01:47)

Lecture-39
Chi-Square Test for Goodness of fit

Suppose we are sampling from a distribution (cdf) $F(x)$ which may depend upon a parameter θ . We want to test

$$H_0: F(x) = F_0(x) \quad \forall x$$
$$H_1: \text{not } H_0$$

where $F_0(x)$ is a known cdf.

We follow the given procedure:

Divide the range of the distribution in k mutually exclusive and exhaustive intervals, say I_1, \dots, I_k .

So, let us have the situation of the form, suppose we are sampling from a distribution. So, let me write distribution function c d f say $F(x)$. Of course, there may be a situation where it depends upon certain parameters. So, we may write that thing, which may depend upon a parameter say θ . So, this θ could be having several components also. So, we want to test, say $H_0: F(x) = F_0(x)$, for all x against H_1 not H_0 . That means, this is not true for some points at least. Where F_0 is a known c d f, that means, we want to test whether the data which has been collected comes from a given distribution $F_0(x)$. So, in the chi square test for goodness of fit the procedure is as follows: We follow the given procedure. So, divide the range of the distribution in k mutually exclusive and exhaustive intervals. Let us call them intervals as I_1, I_2, \dots, I_k .

(Refer Slide Time: 04:31)

Let us assume that $P(X \in I_i) = \pi_i, i=1, \dots, k$.
 Each sample value falls in exactly one of the intervals. Let O_1, \dots, O_k be the respective observed number of observations in the intervals I_1, \dots, I_k .
 Then the vector $\underline{O} = (O_1, \dots, O_k)$ has a multinomial distribution

$$P(O_1 = o_1, \dots, O_k = o_k) = \frac{n!}{\prod_{i=1}^k o_i!} \prod_{i=1}^k (\pi_i^{o_i})$$
 where $\sum o_i = n, \sum \pi_i = 1$
 Also $E(O_i) = n\pi_i = e_i \rightarrow \text{say.}$
 $V(O_i) = n\pi_i(1-\pi_i)$
 $i=1, \dots, k.$

So, now, each value will fall exactly in one of the intervals because, they are mutually exclusive and exhaustive. Let us also assume that, let us assume that probability of x being in an interval I is π_i for i is equal to 1 to k . So, now each sample value falls in exactly one of the intervals. So, let us define the observed frequencies. Let O_1, O_2, \dots, O_k be the respective observed number of observations in the intervals I_1, I_2, \dots, I_k . So, what we are observing? We are observing x_1, x_2, \dots, x_n . Now, you see some of these x_i 's will belong to interval I_1 , some of the x_i 's will belong to interval I_2 and so on. So, we make this break up.

Let O_1 denote the number of x_i 's in the interval I_1 , let O_2 the number of x_i 's in I_2 , let O_k denote the number of x_i 's in I_k . Then, O_1, O_2, O_k these are called observed frequencies of the data. So, now, you see here. You have k categories unlike binomial distributions, where you have 2 categories, here you have k categories. So, if I find out the distribution of O_1, O_2, O_k , this will become a multinomial distribution. So, we will write it in this form, then the vector O , that is O_1, O_2, O_k has a multinomial distribution. n factorial divided by product of O_i factorial, then product of p_i to the power O_i , i is equal to 1 to k where $\sum O_i = n$ $\sum p_i = 1$. That is p_i is the probability of the i th interval. Therefore, probability of O_i is equal to small o_i , that will be given by this multinomial function. Also, we will have from the properties of the multinomial distribution, expectation of O_i , this become equal to $n p_i$, we call it e_i . And variance of O_i , that will become equal to $n p_i (1 - p_i)$, for i is equal to 1 to k .

(Refer Slide Time: 08:22)

Handwritten mathematical derivation on a whiteboard:

$$\text{If } k=2, \frac{O_1 - n\pi_1}{\sqrt{n\pi_1(1-\pi_1)}} \xrightarrow{L} N(0,1)$$

Hence $\frac{(O_1 - n\pi_1)^2}{n\pi_1(1-\pi_1)} \xrightarrow{L} \chi^2_1$

Using $O_2 = n - O_1$ (for $k=2$.)

$$\frac{(O_1 - n\pi_1)^2}{n\pi_1} + \frac{(O_2 - n\pi_2)^2}{n\pi_2} = \frac{(O_1 - n\pi_1)^2}{n\pi_1(1-\pi_1)}$$

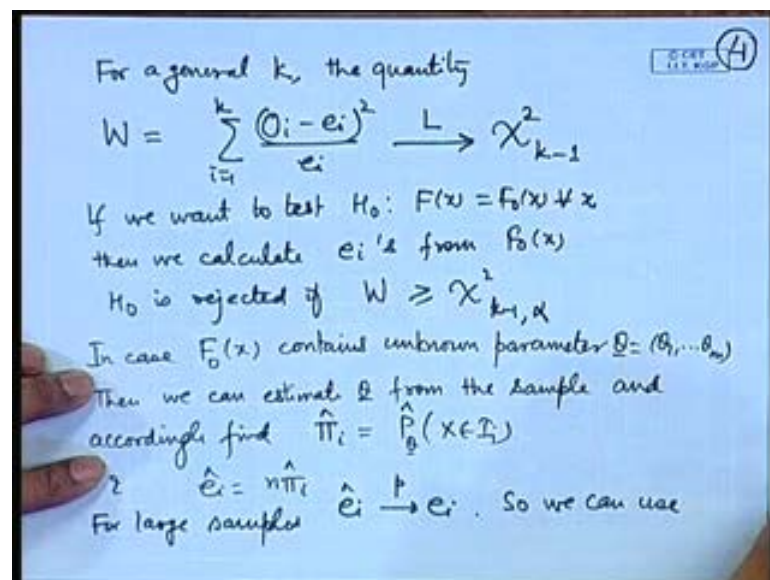
Hence $\sum_{i=1}^2 \frac{(O_i - e_i)^2}{e_i} \xrightarrow{d} \chi^2_1$

Let us take the case of 2 categories. Then, let us see how the test can be conducted. So, if k is equal to 2, then if I look at $x_1 - n p_1$, in place of this I will write O_1 . So, $O_1 - n p_1$ divided by root $n p_1 (1 - p_1)$. This is converging in distribution to normal 0 1. This is from the property of the binomial distribution that the distribution of $x - n p$ divided by root $n p q$ is asymptotically normal. That is the normal approximation to the binomial distribution. So, if we utilize that because for k equal to 2,

this multinomial has become binomial. Therefore, the distribution of O_1 is binomial $n p_1$. So, $O_1 - n p_1$ divided by square root $n p_1 (1 - p_1)$ is asymptotically normal distribution. So, if I take the square of this, then this is asymptotically chi square distribution on 1 degree of freedom, because square of a normal distribution, standard normal variable is a chi square variable.

Now, O_2 is $n - O_1$. This is for the case k equal to 2. So, you can easily see that, if I write down $O_1 - n p_1$ square by $n p_1 (1 - p_1)$ plus $O_2 - n p_2$ square by $n p_2 (1 - p_2)$, then, so, here you substitute O_2 is equal to $n - O_1$ and p_2 is equal to $1 - p_1$. Then, after simplification because, this will become $1 - p_1$. We can take LCM and adjust the terms. This will give simply $O_1 - n p_1$ square divided by $n p_1 (1 - p_1)$. So, what we are observing, that $\sum_{i=1}^k \frac{(O_i - e_i)^2}{e_i}$ and this time I call e_i , i is equal to 1 to k , this is having asymptotically chi square distribution on $k - 1$ degree of freedom.

(Refer Slide Time: 11:03)

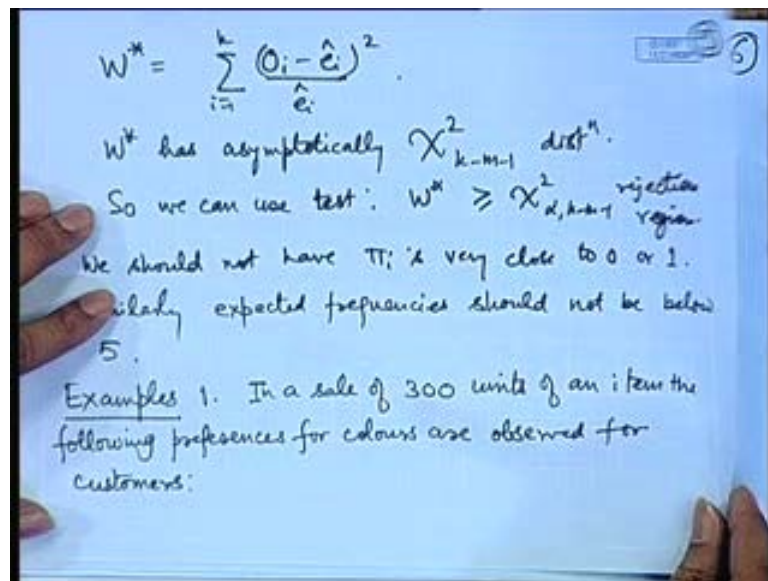


So, if we generalize this thing, in place of 2, if I write for a general k , the quantity, let me call it W . That is equal to $\sum_{i=1}^k \frac{(O_i - e_i)^2}{e_i}$, i is equal to 1 to k . This is having asymptotically chi square distribution on $k - 1$ degrees of freedom. I am using L and d as for the asymptotic distribution here.

So, if we want to test the hypothesis, if we want to test $H_0: F_X$ is equal to F_{naught} , then we calculate e_i 's from F_{naught} . That is, under the distribution F_{naught} what is the probability of the i th interval. So, that is p_i and if I multiply by n , I will get e_i . But, of course, this is calculated **from the** from the known distribution, this one here. So, you can, then consider the difference between the observed frequency and the expected frequency, squared and divided by the expected frequency. So, you can see that, if this hypothesis is true, then the differences between O_i 's and e_i 's must be small. And therefore, this term should be rather small. Therefore, by comparing with the tabulated value of a chi square distribution on $k - 1$ degrees of freedom, we can test whether H_0 can be rejected or cannot be rejected. If it is not true, then these differences will tend to be large. So, the value of W will be large. So, the test H_0 is rejected if W is greater than or equal to $\chi^2_{k-1, \alpha}$.

Now, there may be a situation where F_{naught} may not be completely known. That means, it may include certain parameter. If this includes certain parameter, then from the data we can estimate that parameter also. And then in place of e_i , we can say we are getting e_i heads. And we can substitute there, in case F_{naught} contains unknown parameter, say θ is equal to $\theta_1, \theta_2, \dots, \theta_n$. Then, we can estimate θ from the sample. And accordingly, find p_i head is equal to probability x belonging to I head. That means, the estimate of this. And e_i head is equal to n times p_i head. Then, for large samples, e_i heads will converge to e_i in probability or with probability 1. So, we can use, say W^* that is equal to $\sum (O_i - e_i)^2 / e_i$, i is equal to 1 to k . W^* has asymptotic chi square $k - m - 1$ distribution.

(Refer Slide Time: 15:13)



So, once again, we can use test as that W star greater than or equal to chi square alpha k minus m minus 1. This is the rejection region. However, there are certain precautions one should take while using the chi square approximation. As the binomial approximation to the normal distribution is good, when p is moderate; that means, it should not be close to 0 or close to 1. In a similar way, here the sale probabilities are say π_i 's. So, if either of the π_i 's is extremely small; that means, close to 0 or close to 1, then in that case the expected frequency of that sale will become either too small or too large. If it is too large then for some other sale it may become too small. In that case, this approximation is not good. So, we have the following considerations.

We should not or we can say, we should not have π_i 's very close to 0 or 1. Similarly, expected frequencies should not be below 5. So, a practical consideration that has been done, that there may be a case that firstly, we split the intervals without knowing the probabilities but, when we actually calculate the probabilities and find that the expected frequencies are below 5, then, what we do? We can merge some adjacent intervals. So that, the number of interval becomes slightly less but each sale frequency becomes more than 5. So, this is a practical approach that is used here.

Let me explain this test through certain examples. In a sale of say, 300 units of an item, the following preferences for colors are observed for customers. So, it may be like

certain item of the type, say for example, somebody is buying say car or somebody is buying a say two wheeler. So, we look at the color of the car for example.

(Refer Slide Time: 19:21)

Color preferences of customers.

Brown	Grey	Red	Blue	White	Total
88	65	52	40	55	300

Test the hypothesis that all colours are equally popular.
 Let π_i denote the prob. of i^{th} colour, $i=1, \dots, 5$
 Then $H_0: p_1 = p_2 = \dots = p_5 = \frac{1}{5}$
 $H_1: \text{at least one inequality}$
 $e_i = n\pi_i = 300 \times \frac{1}{5} = 60, i=1, \dots, 5$

$$W = \sum_{i=1}^5 \frac{(O_i - e_i)^2}{e_i} = \frac{28^2}{60} + \frac{5^2}{60} + \frac{(-8)^2}{60} + \frac{(-20)^2}{60} + \frac{(-5)^2}{60}$$

$$= 21.635, \quad \chi^2_{4, 0.05} = 9.487$$

$$\chi^2_{4, 0.01} = 13.28$$
 So H_0 is rejected.
 i.e. customers have colour preference.

So, out of 300 customers, suppose we observe that brown, the colors available are brown grey, red, blue and say white. Out of 300 customers, we find 88 prefer brown, 65 prefer grey, 52 prefer red, 40 prefer blue and 55 prefer white, out of 300. Color preferences of customers. So, we want to test the hypothesis that all colors are equally popular; that means, the customers have equal preference for each of the 5 colors. So, if we want to frame a hypothesis in the form of a test of goodness of fit, what we can do is, the hypothesis is of the form that each sale has probability 0.2.

So, let p_i denote the probability of i^{th} color, for i is equal to 1 to 5. There are 5 colors here. Then, we want to test that each of the, is 1 by 5. Now, based on these assumption, so, here basically the sale or interval is actually the type here. So, brown is one type, grey is another type, red is another type, blue is another type, white is another type. So, this is again a multinomial situation and here we are assuming the probabilities to be same in the null hypothesis. H_1 is that, at least 1 inequality. So, on the basis of this, we do the following calculations. We can calculate e_i 's. So, e_i is $n p_i$. So, here it is 300 into 1 by 5, that is equal to 60. Each category has the same probability and therefore, each category will have the same expected frequency also.

So, on the basis of this if we calculate W , that is $\sum (O_i - e_i)^2 / e_i$, it is equal to 1 to 5. Then, this is equal to $88 - 60$, that is $28^2 / 60$ plus $65 - 60$, second sale, that is $5^2 / 60$, third sale is 52 . So, it is $8^2 / 60$ plus $40 - 60$ is -20 , square by 60 plus $55 - 60$ is -5 , square by 60 . So, the sum can be easily evaluated, it turns out to be 21.635.

Now, there are 5 categories. We also notice that the expected frequency of each sale is more than 5. So, the chi square assumption is valid. Therefore, we look at the value of chi square on 4 degrees of freedom. Suppose, I consider the value at say 0.05, then from the tables of the chi square distribution, one can find this value is 9.487. We may even look at say chi square 4.01, that is 13.28. So, you can see that the calculated value of W , that is 21.635 is bigger than this. So, H_0 is rejected; that means, what is the conclusion? The conclusion is that customers have preferences for the colors.

You can see the raw data here. The observed frequency for brown is 88, which is almost more than twice the choice of blue color. If we see the choice of grey, that is much higher. The choice of red, blue and white is below. So, you can see that, in specifically speaking, brown and blue, they cause major discrepancies here. Blue is, say least favorable color and brown is the most favorable color here. In fact, if I have only 3 of these, then they look almost nearby 65, 52 and 55. That is, customers have color preferences. So, basically what we have tested is something like a discrete uniform distribution. And we conclude that the data does not follow a discrete uniform distribution.

(Refer Slide Time: 25:17)

2. For a particular organism, three types of genotypes A, B, C are possible. A theory suggests that they may be in the ratio 1:2:1. To test this theory a sample of 90 units are taken with the following results! Test whether the data support the theory.

	Genotypes			Total
	A	B	C	
O_i	18	44	28	90
π_i	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{1}{4}$	90
e_i	22.5	45	22.5	90

So $H_0: p_A = \frac{1}{4}$
 $p_B = \frac{2}{4}$
 $p_C = \frac{1}{4}$

Can't be rejected
 i.e. data support the theory.

$\sum_{i=1}^3 \frac{(O_i - e_i)^2}{e_i} = 2.26$, $\chi^2_{2, 0.05} = 5.99$,
 $\chi^2_{2, 0.01} >$

Let us take another example, for a particular organism three types of genotypes A, B and C are possible. A theory suggests that, they may be in the ratio, say 1 is to 2 is to 1. Now, to test this hypothesis, to test this theory, a sample of 90 units **are** is taken with the following results. So, genotypes A, B, C, it is observed that out of 90 units, 18 had genotype A, 44 had genotype B, and 28 had genotype C. The total is 90. So, now, we want to test whether the data supports the theory. So, for this once again, we have 3 categories. The probabilities of each category, let me call it π_i . So, 1 is to 2 is to 1. So, this probability is 1 by 4, this probability is 2 by 4, that is half and this probability is 1 by 4. So, expected frequency, this is observed frequency. We can actually do the calculations in the form of a table here.

So, here you see, if the probability of the genotype A is 1 by 4, the total number of units is 90. So, the expected frequency for that will be 90 by 4. That is 22.5 here, it will become 45, here it will become 22.5. That is, total is 90. So, based on this, one can carry out the calculations, $\sum O_i - e_i^2$ by e_i , i is equal to 1 to 3. So, for example, here it will become 4.5 square divided by 22.5 plus 1, 44 minus 45. So, 1 square by 45 and 28 by minus 22.5, that is 5.5 square by 22.5. So, one can look at this calculations, this turns out to be 2.26.

So, one can easily compare with the chi square value on 2 degrees of freedom. Suppose, we look at 0.05, then this is 5.99. And of course, if I look at say chi square 2 at 0.01, this is going to be larger than this. So, we cannot reject H_0 ; that means, H_0 that is p A is equal to 1 by 4, p B is equal to half, p C is equal to 1 by 4, cannot be rejected. That means, the data supports the theory that the genotype are in the proportion 1 is to 2 is to 1.

(Refer Slide Time: 30:06)

The whiteboard shows the following derivation:

$$W = \sum_{i=1}^k \frac{(O_i - e_i)^2}{e_i} = \sum \left\{ \frac{O_i^2 + e_i^2 - 2O_i e_i}{e_i} \right\}$$

$$= \sum \frac{O_i^2}{e_i} + \sum e_i - 2 \sum O_i$$

$\sum e_i = n$
 $\sum O_i = n$

$$= \sum_{i=1}^k \frac{O_i^2}{e_i} - n \quad (k-1 \text{ d.f.})$$

$$W^* = \sum_{i=1}^k \frac{(O_i - \hat{e}_i)^2}{\hat{e}_i} = \sum_{i=1}^k \frac{O_i^2}{\hat{e}_i} - n \quad (k-m-1 \text{ d.f. of unknown parameters are there})$$

Here, one point about this calculation also. This formula which we have given here, $\sum (O_i - e_i)^2 / e_i$, one can actually have an alternative form for this. We can consider expanding this term. So, it is $O_i^2 / e_i + e_i - 2O_i$ and then summation. So, summation e_i is actually n , minus twice summation O_i . So, summation e_i is n and summation O_i is also n . So, this becomes simply $\sum (O_i^2 / e_i) - n$, where i is equal to 1 to k . So, this is an alternative formula for W^* , for W . Similarly, if I am considering W^* , that is $\sum (O_i - \hat{e}_i)^2 / \hat{e}_i$, then once again, this can also be written as $\sum (O_i^2 / \hat{e}_i) - n$. Here, degrees of freedom are $k - m - 1$ and here the degrees of freedom are $k - m - 1$, if m unknown parameters are there. Many times, this expression is easier to calculate.

(Refer Slide Time: 31:58)

Example: One wants to investigate the distⁿ of the number of claims for medical treatments by families. A previous study suggested that the distⁿ may be Poisson. A random sample of 200 families is taken with the following classification

	0	1	2	3	4	5	6	7	Total
	22	53	58	39	20	5	2	1	200

We want to test whether a Poisson distⁿ fits the data appropriately.

First we estimate the parameter λ of the Poisson distⁿ. MLE for λ is $\bar{x} = \frac{0 \times 22 + 1 \times 53 + \dots + 7 \times 1}{200} = 2.05 \approx 2.0$

Let us take one case, where the distribution will depend upon certain unknown parameter. So, one wants to investigate the distribution of the number of claims, for medical treatments by families. A previous study suggested that the distribution may be poisson. So, to investigate this, a random sample of 200 families is taken with the following classification. So, that is, for each family how many claims are there? And that is what is the frequency? How many families made how many claims? So, it turned out that 22 families did not make any claim, 53 families made 1 claim, 58 families made 2 claims, 39 families made 3 claims, 20 families made 4 claims, 5 families made 5 claims, 2 families made 6 claims, and 1 family made 7 claims; that means, no family made more than 7 claims over a period. Over a given period, may be say 5 years time or may be over a 10, 10, etcetera.

So here, we want to test whether the, we want to test whether a poisson distribution fits the data appropriately. Now, here we do not even know the parameter lambda of the poisson distribution. So firstly, we will estimate that. So, first we estimate the parameter lambda of the poisson distribution. Now, we have already done estimation. We may use say maximum likelihood estimator or say minimum variance unbiased estimator for lambda. So, for example, maximum likelihood estimator for lambda is \bar{x} . So, \bar{x} is the mean which can be evaluated from here. That is 0 into 22 plus 1 into 53 and so on, plus 7 into 1 divided by 200. The total number of families is 200. So, this value turns out

to be 2.05. So, we may approximately take 2 as the lambda value and we like to check whether this data follows poisson at the rate 2. Now, this is a reasonable approximation because, 2.05 is a value and here, we are talking about the number of claims. So, it is appropriate to take lambda to be an integral approximation of the value and it is extremely close. So, this is fine.

(Refer Slide Time: 36:51)

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

	0	1	2	3	4	5	6	7 & above	Total
\hat{p}_i	0.135	0.271	0.271	0.18	0.09	0.036	0.012	0.005	1
\hat{e}_i	27	54.2	54.2	36	18	7.2	2.4	1.0	200
							10.6		

$$W^* = \sum_{i=1}^6 \frac{(O_i - \hat{e}_i)^2}{\hat{e}_i} = 2.33$$

$$\chi_{(6-1)}^2 = \chi_{4, 0.05}^2 = 9.487$$

So H_0 cannot be rejected
 i.e. Poisson distribution seems to be an appropriate model for the data.

So, we want to find out the probabilities of x is equal to, say small x for 0, 1, etcetera. So, now, the formula in the poisson distribution is e to the power minus lambda, lambda to the power x by x factorial, x is equal to 0, 1, 2 and so on. Now, the point here is that, this is an infinite value distribution. Here, we will calculate the probabilities only for 0, 1 to 7. So, they will not add up to 1. So, what we will do, we can calculate the probabilities 0 1 to upto 6, and 7 we will put into the form of 7 and above. Although, the data is not collected like that, it was observed that for 8, there is no family which made 8 claims, there is no family which made 9 claims, etcetera. So, basically the values 7 corresponds to 7 and above. In that case, the probability of the events, each of the classes will be adding up to 1. So, we calculate the probabilities and we report it like this. So, p i's for 0 1, 2, 3, 4, 5, 6, 7 and above.

So, substituting the value of lambda as lambda head is equal to 2, we can calculate, this is e to the power minus 2, that is equal to 0.135. The probability of x is equal to 1, that

will become $\lambda e^{-\lambda}$ to the power minus λ . So, that is twice $e^{-\lambda}$ to the power minus λ . So, it is 0.27, this value is 0.271 again, because this is again $2 e^{-\lambda}$ to the power minus 2, this is 1.8, 0.09, then 0.036, 0.012, 0.005. The total is 1 here. So, now the expected frequency is, the estimate of the expected frequencies can be calculated by multiplying by 200 here. So, we will get it as 27, 54.2, 36, 18, 7.2, 2.4, and 1.0. The total is 200. Now, you note here. The sale frequency of sixth and seventh, 6 and 7 and above, these are much below 5. In fact, even if we merge these 2 sales, the expected frequency will be only 3.4. So, then this will not allow us to use the chi square approximation here. So, what we do, we merge the last 3 sales. So, if we merge these last 3 sales and add up, we will get 10.6 as the expected frequency, which is above 5. So, in place of 8 sales, now we have only 6 sales. So, $\sum (O_i - E_i)$. So, we will have to then merge the observed frequency also here. The observed frequency merged will give 5 plus 2 plus 1, that is equal to 8. So, now, $\sum \frac{(O_i - E_i)^2}{E_i}$, 1 to 6 only. This is W^* . That is equal to 2.33.

Now, chi square value has to be on $k - m - 1$. So, here minus 1 because only 1 parameter is there k is 6. So, this becomes 6 minus 1 minus 1 that is chi square 4. Now, if we see the value at say 0.05, that is 9.48, etcetera. So, H_0 cannot be rejected; that is poisson distribution seems to be an appropriate fit or appropriate model for the data. At least on the basis of the given data, we have no reason to reject H_0 that, the distribution comes from a poisson.

(Refer Slide Time: 42:40)

$r \times c$ Contingency Tables

Sometimes the data in a statistical expt is categorical rather than numerical.

$r \rightarrow$ rows, c columns

	1	2	...	c
1	O_{11}	O_{12}		O_{1c}
2	O_{21}	O_{22}		O_{2c}
...				
r	O_{r1}	O_{r2}		O_{rc}

The sample n is categorized in rc cells.
 $O_{ij} \rightarrow$ observed frequency of the $(i,j)^{th}$ cell.

investment options

The chi square test for goodness of fit has other applications also. Let us consider the data which is not quantitative, but rather qualitative in nature. Let us consider the situation of contingency tables. r by c contingency tables. Sometimes the data in a statistical experiment is categorical rather than numerical. For example, in a population of individuals, we will like to know how many of the people are smokers and how many of them are non-smokers. So, it is a categorical data; that is, there is two categories of person, one who are smokers, one who are non-smokers. We may also categorize them according to those, who are, who ultimately get lung cancer and who do not get the lung cancer. So, now, the situation is like this in a population, we have characterized according to two different methods of categorization. One is the smoking habit and another is the incidence of a disease.

Now, we want to check whether there is any association between these two characteristics or two attributes; that means, is it true or is it found on the basis of the data, that those who smoke, they get, they are more likely to get a lung cancer? So, this is called testing for independence in a contingency table. So, what is a contingency table? That we have two types of categorizations. One is say, we will represent in the form of rows and another type in the category, in the columns. So, we say r rows and c columns. So, the data may be represented like this. So, we have category A, we have category B. This may be 1, 2 up to c and here, you will have 1, 2 up to r .

For example, in an organization people are classified according to their professional hierarchical levels. For example, one may be the top person of the organization, that is say top management, in 2 we may have some people who are in the supervisory position or executives and here, you may have say daily wage workers. And now, we want to look at the attitude towards a certain option. So, for example, there is a option given for them for investment, investment option; that means, from their salary, they may be allowed to invest in certain, say shares of the company. And there may be, say 4 different types of shares. One is where you get a fixed kind of return, another is where you have a safe return; that means, the equity is distributed into safe instruments, next may be some balanced; that is partially risk and partially balanced, partially safe. And there may be risky option. So, maybe you will find that, out of a sample of n implies of the organization, you found O_{11} are of, that is those who are in the top management and those who go for option 1, O_{12} , O_{1c} , O_{21} , O_{22} , O_{2c} , O_{r1} , O_{r2} , O_{rc} . So, this is the data. That is, the sample n is categorized in $r \times c$ sales and O_{ij} is the observed frequency of the i jth sale.

(Refer Slide Time: 48:12)

We want to test whether there is any association in two ways of categorization (A & B) (12)

Testing for independence in a contingency table

$\pi_{ij} \rightarrow$ exp prob. of $(i, j)^{th}$ cell

$\pi_{i.} \rightarrow$ prob of i^{th} type of Cat B

$\pi_{.j} \rightarrow$ prob of j^{th} type of Cat A.

$\pi_{i.} = \sum_{j=1}^c \pi_{ij}$, $\pi_{.j} = \sum_{i=1}^r \pi_{ij}$

Then we want to test $H_0: \pi_{ij} = \pi_{i.} \cdot \pi_{.j}$ for every pair (i, j)

$H_1:$ at least one inequality.

What we want to test here is, we want to test whether there is any association in two ways of categorization, that is A and B. Whether A and B are having any association. As I just gave the example of smoking and cancer, similarly here may be that, those who are

in the higher order in the hierarchy of the company, they may like to invest in high risk equities, etcetera. Whereas, those who are in the lower income group or lower in the hierarchy, they may like to go for safe, this one. Suppose, this could be our hypothesis or we may make a hypothesis that there is no relation.

So, this is called testing for independence in a contingency table. So, you can say here, p_{ij} is the expected or you can say probability of i j th sale; that is the theoretical probability. $p_{i \cdot}$ is the probability of i th type of category, say B. $p_{\cdot j}$ is the probability of j th type of category A. So, actually here, $p_{i \cdot}$ will be equal to $\sum p_{ij}$, $\sum_j p_{ij}$ is equal to 1 to c and $p_{\cdot j}$ is actually $\sum p_{ij}$, $\sum_i p_{ij}$ is equal to 1 to r . Then, we want to test $H_0: p_{ij} = p_{i \cdot} p_{\cdot j}$ for every pair i, j and H_1 is at least 1 inequality.

Now, if we see carefully, this is nothing but a generalization of the test of goodness of fit itself. In the test of goodness of fit, what we wanted to test is that whether the data comes from a particular distribution. The procedure adopted was that, we divided the range of the distribution into k intervals or you can say k sales, and we looked at the observed frequencies. The theoretical distribution of the observed frequencies was multinomial. Now, likewise here, if you see, what we are claiming here is that the probabilities are p_{ij} 's, for the i, j th sale the observed frequency is O_{ij} . So, if we look at the distribution of the categories here, that is O_{11}, O_{12}, O_{1c} and so on upto O_{rc} . The joint distribution of this will again be multinomial. So, that means, the test will be actually of the previous form itself with little bit modification.

(Refer Slide Time: 52:23)

$(O_{11}, O_{12}, \dots, O_{rc})$ has a multinomial distⁿ with cell probabilities π_{ij} , $\sum \sum \pi_{ij} = 1$.
 $e_{ij} = n \pi_{ij}$.
 Then $\tilde{W} = \sum \sum \frac{(O_{ij} - e_{ij})^2}{e_{ij}} \longrightarrow \chi^2_{(r-1)(c-1)}$
 $\hat{e}_{ij} = \frac{R_i C_j}{n}$, $R_i = \sum_{j=1}^c O_{ij}$, $C_j = \sum_{i=1}^r O_{ij}$
 $\tilde{W} = \sum \sum \frac{(O_{ij} - \hat{e}_{ij})^2}{\hat{e}_{ij}} \longrightarrow \chi^2_{(r-1)(c-1)}$
 So: $R_i C_j H_0 \nabla \tilde{W} \geq \chi^2_{(r-1)(c-1), \alpha}$.

So, if we write down O_{11} , O_{12} , and so on upto O_{rc} . This has a multinomial distribution with cell probabilities π_{ij} with $\sum \pi_{ij} = 1$. So, using the previous argument, if I define say e_{ij} is equal to n times π_{ij} , then double summation $\sum \sum \frac{(O_{ij} - e_{ij})^2}{e_{ij}}$; that means, call it say \tilde{W} . Then, this is asymptotically chi square distribution on $(r-1)(c-1)$ degrees of freedom.

The values of e_{ij} 's are estimated by considering $R_i C_j$ by n , where R_i is actually summation of O_{ij} , j is equal to 1 to c and C_j is $\sum O_{ij}$, i is equal to 1 to r . So, if we make use of this, $\sum \sum \frac{(O_{ij} - \hat{e}_{ij})^2}{\hat{e}_{ij}}$, then this will be asymptotically chi square on $(r-1)(c-1)$ degrees of freedom. So, the test is reject H_0 , if \tilde{W} is greater than or equal to $\chi^2_{(r-1)(c-1), \alpha}$.

So, the chi square test for goodness of fit is applicable for testing for independence in a r by c contingency table also. So, in the next class I will give some applications of this test. We will also see that, this is applicable to a slightly different situation. Here, you can see that we are considering a random sample of size n from the population and then we are making this classification. In certain other situations, even this row totals or column totals may be fixed. And then the sampling is done from there; that means, something like a stratified sampling, we will see that even in the stratified sampling the same formula is applicable. So, in the next class I will be covering those portions.