

**Statistical Inference**  
**Prof. Somesh Kumar.**  
**Department of Mathematics**  
**Indian Institute of Technology, Kharagpur**

**Lecture – 60**  
**Test for Goodness of Fit – II**

If we considered the large sample test actually it is a chi square test. Let me give this thing.

(Refer Slide Time: 00:26)

Testing for Independence in  $r \times c$  Contingency Tables

A	B	$B_1$	$\dots$	$B_c$	$R_i$
$A_1$	$O_{11}$			$O_{1c}$	
$\vdots$					
$A_r$	$O_{r1}$			$O_{rc}$	$R_r$
	$C_1$			$C_c$	$N$

$O_{ij} \rightarrow$  observed frequency of the  $(i,j)^{th}$  cell  
A r. s. of size  $n$  is taken from the total pop'n.

Suppose a population is categorized according to two attributes A and B.  
eg A  $\rightarrow$  Income level      A  $\rightarrow$  two four levels  
 $A_1 \rightarrow$  high income       $A_2 \rightarrow$   $\dots$   
 $A_2 \rightarrow$  Upper Mid       $A_3 \rightarrow$   $\dots$   
 $A_3 \rightarrow$  Lower       $A_4 \rightarrow$   $\dots$

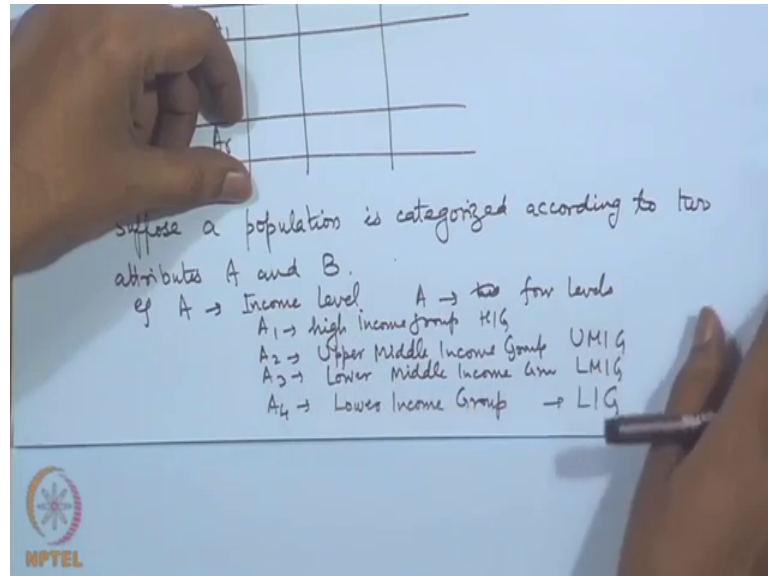
UMIG  
LMIG

So, testing for independence in  $r$  by  $c$  contingency tables. So, contingency tables are used when the data is categorical. And so, for example, you may like to test whether the selection of say students in a particular entrance examination dependent upon their financial status, or you can say the financial status of their parents, or is it dependent on the educational status of their parents or is it dependent upon the level of the schools where they are studying the school board and so on.

So, on one side you have certain category let us say  $A_1 A_2 A_r$ . On the other side you have category B which is having say  $c$  classifications. So, this is attribute A this is attribute B. Suppose a population is categorized according to 2 attributes A and B. Say for example, A is the income level. So, we may say A has say 4 levels. Say  $A_1$  is High Income Group say HIG.  $A_2$  is say Upper Middle Income Group that is say UMIG then a

3 say Lower Middle Income Group L M I G and say A 4 is lower income group that is say L I G.

(Refer Slide Time: 03:05)



So, I am distributing or you can say dividing the population into 4 groups according to the attribute which is income level.

(Refer Slide Time: 03:27)

B → Expenditure on Education

- B<sub>1</sub> → More than normal
- B<sub>2</sub> → Average
- B<sub>3</sub> → below average

We want to test whether there is a relation (association) between the two attributes.

A \ B	B <sub>1</sub>	B <sub>2</sub>	B <sub>3</sub>	Total
A <sub>1</sub>	78	80	85	243
A <sub>2</sub>	112	110	105	327
A <sub>3</sub>	110	100	95	305
A <sub>4</sub>	30	50	45	125
Total	330	340	330	1000

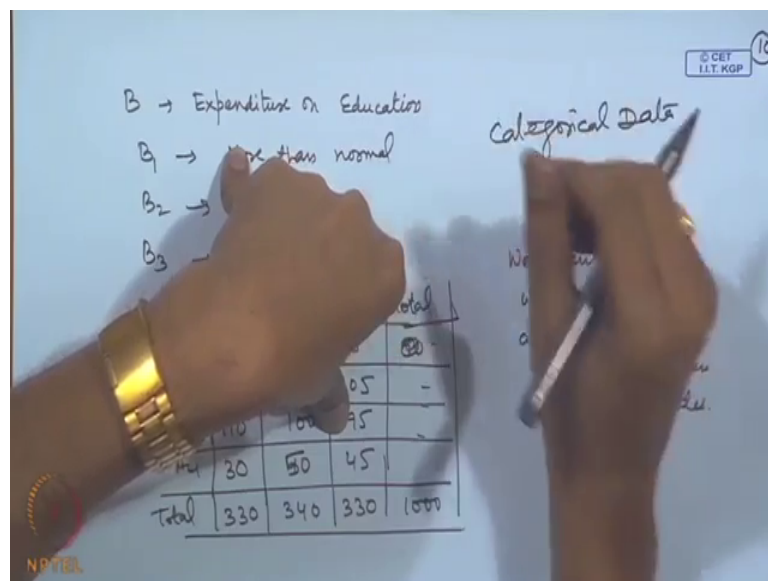
And second one is say expenditure on say education. And then again I put B 1 B 2 B 3. So, this is say more than normal average and below average I consider. So, this is A this is B now A 1 A 2 A 3 A 4. On this side you have B 1 B 2 B 3. I consider say a random

sample of say thousand individuals. Out of the thousand individuals how many fall into each category. How many fall into each category. So, there are say 12 classes here. I am just creating an hypothetical data here say 78 1 1 2.

So, this is say let me put 3 classes are there. Let me put some rough number here say 330 340 and say 330. This is just incidental I am fixing the values in some way so, that the total match here. So, this is 190 plus 1 1 0 say 30 80 1 1 0 and say 140 say this is 85, 1 0 5 95. So, this is 140 this is 50 here this is say 90 sorry. So, this totals will add up here.

Now, we want to test whether there is a relation or you can say association between the 2 attributes. Can you say that if the income is high the expenditure on the education will be more? Or if the income is low then the expenditure on the education will be less and so on; that means, is there a dependence on this. So, we want to test this. You can see this is called a categorical data.

(Refer Slide Time: 06:50)



Categorical data as oppose to the numerical data that you have in the other problems.

Now, for this we can create the situation like this. So, I am making a general setup. You have observed frequencies I will call  $O_{11}$   $O_{12}$   $O_{21}$   $O_{22}$  these are the  $O_{ij}$  is the observed frequency of the  $ij$  th cell. So, this is  $n$  here. Now I have considered a random sample of size  $n$ . A random sample of size  $n$  is taken from the population ok.

(Refer Slide Time: 07:57)

Let  $\pi_{ij} = P(X \in (i,j)^{\text{th}} \text{ cell}) = P(X \in A_i \cap B_j)$

$H_0$ : Row & col attributes are independent

$H_1$ :  $\rightarrow$  not

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

$\pi_{i.} \rightarrow P(X \in A_i)$   $\pi_{.j} = P(X \in B_j)$

So the hypothesis of independence is then equivalent

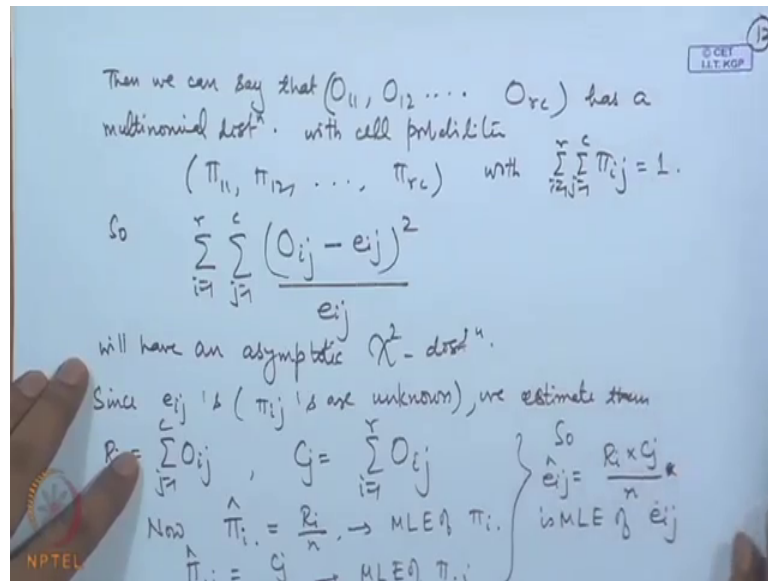
$\Rightarrow \pi_{ij} = \pi_{i.} \times \pi_{.j} \quad \forall i, j$

Now the expected frequency of  $(i,j)^{\text{th}}$  cell  $= n \pi_{ij} = e_{ij}$

Let us consider say let  $\pi_{ij}$  be the probability of the observation belonging to  $i$   $j$ th cell. Our null hypothesis is that row and column attributes are independent. And our alternative hypothesis is that they are not independent.

So, we write down the totals here  $\pi_{i.}$  that is equal to  $\sum_{j=1}^c \pi_{ij}$  is equal to  $1/2/c$   $\pi_{.j}$  that is equal  $\sum_{i=1}^r \pi_{ij}$  over  $I$  is equal to  $1/r$ . So,  $\pi_{i.}$  is the probability of observation belonging to the  $A_i$ . And similarly  $\pi_{.j}$  this is the probability of  $X$  belonging  $B_j$  ok. And  $\pi_{ij}$  is probability of  $X$  belonging to  $A_i$  intersection  $B_j$  kind of thing. So, the hypothesis of independence is then equivalent to that  $\pi_{ij}$  is equal to  $\pi_{i.}$  into  $\pi_{.j}$  for all  $\pi_{ij}$ . Now the if independence is there expected frequency of  $i$   $j$ th cell that will be equal to  $n$  times  $\pi_{ij}$  let me call it  $e_{ij}$  ok.

(Refer Slide Time: 10:17)



So, then we can say that  $O_{11}, O_{12}$  and so on all of this taken together this has a multinomial distribution with cell probabilities  $\pi_{11}, \pi_{12}$  and so on  $\pi_{rc}$ . And sigma of  $\pi_{ij}$  is equal to 1. So, by the goodness of fit explanation that I gave earlier, if I considered double summation  $\sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - e_{ij})^2}{e_{ij}}$  will have an asymptotic chi square distribution. Now how do you evaluate these things? Since the theoretical values are unknown, that is  $e_{ij}$  that is  $\pi_{ij}$  are unknown we estimate them ok.

How will you estimate that? You can consider say  $R_i$  that is equal to  $\sum_{j=1}^c O_{ij}$ , that is the totals on the row side. If I add here I call it  $R_1$  and so on  $R_r$  and on this side if I add I call it  $C_1, C_2, \dots, C_c$ . That is a row totals and column totals  $\sum_{j=1}^c O_{ij}$  is equal to one to  $C$  and  $\sum_{i=1}^r O_{ij}$  is equal to one to  $r$ . Now  $\hat{\pi}_{i \cdot}$  that is the estimate of the  $i$ th row probability row attribute probability that can be estimated by  $R_i/n$ . This is the maximum likelihood estimator of  $\pi_{i \cdot}$ . Similarly, if I consider  $\hat{\pi}_{\cdot j}$  that is  $C_j/n$  ok.

(Refer Slide Time: 13:10)

multinomial dist<sup>n</sup>. with cell probabilities  $(\pi_{11}, \pi_{12}, \dots, \pi_{rc})$  has a

$(\pi_{11}, \pi_{12}, \dots, \pi_{rc})$  with  $\sum_{i=1}^r \sum_{j=1}^c \pi_{ij} = 1$ .


So 
$$\sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - e_{ij})^2}{e_{ij}}$$

will have an asymptotic  $\chi^2$ -dist<sup>n</sup>. (on  $(rc-1)$  d.f.)

Since  $e_{ij}$ 's ( $\pi_{ij}$ 's are unknown), we estimate them

$R_i = \sum_{j=1}^c O_{ij}$ ,  $C_j = \sum_{i=1}^r O_{ij}$

Now  $\hat{\pi}_{i.} = \frac{R_i}{n} \rightarrow$  MLE of  $\pi_{i.}$  } So  $\hat{e}_{ij} = \frac{R_i \times C_j}{n}$  }  
 $\hat{\pi}_{.j} = \frac{C_j}{n} \rightarrow$  MLE of  $\pi_{.j}$  } is MLE of  $e_{ij}$  }  
 $\left. \begin{array}{l} \sum \pi_{i.} = 1 \\ \sum \pi_{.j} = 1 \end{array} \right\}$



This is the maximum likelihood estimator of  $\pi_{i.}$

So, we can consider  $\hat{e}_{ij}$  as equal to  $R_i$  into  $C_j$  divided by  $n$  into. So, here divided by  $n$  and here divided by  $n$  and then multiplied by  $n$ . So, that  $n$  will cancel out I can write just thing. So, this is MLE of  $e_{ij}$ . So, how many parameters we are estimating here. This total is having  $rc$  classes. So, we will have the degrees of freedom  $rc$  minus 1 this will have on  $rc$  minus 1 degrees of freedom, but now we have estimated  $r$  plus  $c$  terms.


(Refer Slide Time: 14:17)

Since we have estimate  $r+c-2$  parameters, the dof for  $\chi^2$  will be  $rc - 1 - r - c + 2 = (r-1)(c-1)$ .

So 
$$W^* = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - \hat{e}_{ij})^2}{\hat{e}_{ij}} \rightarrow \chi^2_{(r-1)(c-1)}$$

So the test is to reject hypothesis of independence

$\Rightarrow W^* > \chi^2_{(r-1)(c-1), \alpha}$



Since we have estimated  $r + c$  parameters, the degrees of freedom for chi square will be  $rc - 1 - r - c$ . That is equal to  $(r - 1)(c - 1)$  here. I think I have made one error here this will become  $rc - 1$  minus this one.

So, asymptotic distribution of say let me call it  $\chi^2$  that is equal to double summation  $\sum_{i,j} (o_{ij} - e_{ij})^2 / e_{ij}$  where  $\sum_i o_{ij} = 1$  to  $r$  and  $\sum_j o_{ij} = 1$  to  $c$ . That will be asymptotically chi square on  $(r - 1)(c - 1)$  degrees of freedom. This is actually  $(r - 1)(c - 1)$  here. The reason is that the total is going to be 1. We are getting  $\sum_i \pi_i = 1$  and  $\sum_j \pi_j = 1$ . So, the last 2 parameters are not estimated ok. So, we have estimated  $r + c - 2$ . So, this is  $rc - 1 - r - c + 2$  that is  $(r - 1)(c - 1)$  here.

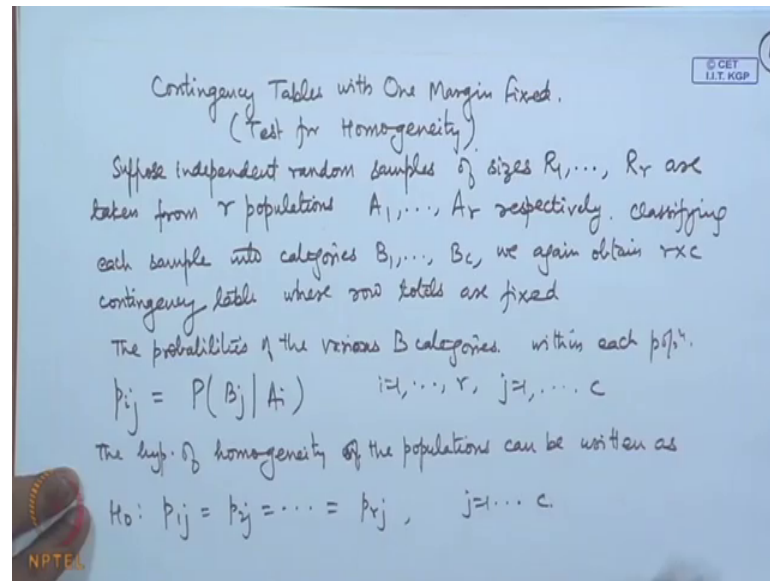
So, asymptotically this will have this distribution. So, the test is to reject hypothesis of independence. If this  $\chi^2$  is greater than  $\chi^2_{(r-1)(c-1), \alpha}$ . Now there is yet either situation here. Because here in this particular case what I have considered the total sample size for the population was fix like  $n$ . Therefore, each of this  $o_{ij}$  these became random variables. These are also random variables and therefore, they had multinomial distributions.

But there can be a situation for example, the population may be large therefore, we may not fix the total sample. Rather we may fix the sample according to the categories. That is called testing for the homogeneity in a contingency table. It could be like this for example, we may like to we have a distribution of people having different diseases ok.

Now, we categorize the population according to different ethnic groups. So, for example, we are considering say Asians, we may consider say Caucasians we may consider say you say Africans and so on. Now we look at that whether in those populations the distribution of the diseases various diseases is the same or not, but in place of taking the random sample of from the total population, we consider a certain number from say Asians, certain number from Africans, certain number from Caucasians and so on. And then will look at the classification among that. Therefore, now the  $\pi_i$  are cgs will not be random, but the individual entries will become random.

So, the setup will be slightly modified; however, we will show that the test statistic is still same.

(Refer Slide Time: 18:36)



So, contingency tables with one margin fixed this is called test for homogeneity. So, here we are considering suppose independent random samples of sizes  $R_1, R_2, \dots, R_r$  they are taken from  $r$  populations say  $A_1, A_2, \dots, A_r$  respectively. In the earlier case I have considered the full population. And then the categories we are called  $A_1, A_2, \dots, A_r$  and according to the other attribute  $B_1, B_2, \dots, B_c$ . Now here from each of them I am considering the samples. So, this will be called populations and after the sample is taken this is classified into  $B_1, B_2, \dots, B_c$  that is according to the attribute  $B$ . So, classifying each sample into categories say  $B_1, B_2, \dots, B_c$  we again obtain  $r$  by  $c$  contingency table.

So, the nature is the same; however, the interpretation of the values is different. Where row totals are fixed. Alternatively, you may fix the column totals and then definitely; that means,  $B_1, B_2, \dots, B_c$  may have fixed sample sizes. And then you can classify each of them according to  $A_1, A_2, \dots, A_r$  ok. Now here the probability then will become conditional probability.

So, probability of the various  $B$  categories within each population they are now  $p_{ij}$  that we call probability of  $B_j$  given  $A_i$  for  $i$  is equal to 1 to  $r$ ,  $j$  is equal to 1 to  $c$ . And the hypothesis of homogeneity of the populations, this can be written as say  $p_{1j}$  is equal to  $p_{2j}$  and so on is equal to  $p_{rj}$  for  $j$  is equal to 1 to  $c$ .



(Refer Slide Time: 22:14)

$$\hat{p}_{1j} = \hat{p}_{2j} = \dots = \hat{p}_{rj} = \frac{c_j}{n} \quad \text{under } H_0.$$

$$\hat{E}_{ij} = R_i \times \hat{p}_{ij} = \frac{R_i \times C_j}{n}$$

$$W^* = \sum \sum \frac{(O_{ij} - \hat{E}_{ij})^2}{\hat{E}_{ij}}$$

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

So this test  $\chi^2$  is the same as the previous one.

Now, for estimating these probabilities we can consider  $\hat{p}_{1j}$  is equal to  $\hat{p}_{2j}$  is equal to  $\hat{p}_{rj}$  that is equal to  $c_j$  by  $n$ . This is under  $H_0$  we are considering. So, for  $i$ th cell  $\hat{E}_{ij}$  this will become  $R_i$  into  $\hat{p}_{ij}$  that is equal to  $R_i$  into  $C_j$  by  $n$ . So, it becomes the same value which I obtained in the previous case. And therefore, the test statistic is same as  $W^*$  that is double summation  $O_{ij}$  minus  $\hat{E}_{ij}$  square divided by  $\hat{E}_{ij}$ . The calculation of the degrees of freedom is slightly different now. You are having  $c$  cells and therefore, in each cell you will have each this one you will have  $c$  minus 1 degrees of freedom.

So,  $r$  times that. So, total degrees of freedom will become  $rc$  minus 1. Now in each of them you will be considering the estimated number of parameters that is  $c$  minus 1. So, that is becoming  $r$  minus 1  $c$  minus 1. So, it is the same degrees of freedom. So, this test statistic is the same as the previous one. So, what we are observing that testing for the independence in a  $r$  by  $c$  contingency table we are getting asymptotic chi square test. And the test is similar when the interpretation of the or you can say the sampling scheme is slightly modified. In place of taking the sample from the full population I distributed into the stratum and then in each strata I take a fixed sample size.

However, the testing procedure does not get modified. When I discuss the various problems on the testing of hypothesis I will discuss the applications of these 2 procedures which I have described today. In the following lecture I will be discussing

another important test that is called a sequential probability ratio tests. So, that I will be covering in the next lecture. I will divert of you lectures and on the problems also for various tests that we have derived in this particular course.