

Statistical Methods for Scientists and Engineers
Prof. Somesh Kumar
Department of Mathematics
Indian Institute of Technology – Kharagpur

Lecture - 15
Parametric Methods - VII

We have discussed in detail the tests for the parameters of normal population. I considered one sample problem, in which we considered the testing for the mean and variance of one normal population. We also considered 2 normal populations and we considered various tests for comparing the means and also the variances. However, when we have qualitative data, we may also be interested in testing for the proportions.

(Refer Slide Time: 00:51)

Lecture - 15

Tests for Proportions : $X \sim \text{Bin}(n, p)$
where n is known. When n is small, then we can consider test based on X .

Define $P = \frac{X}{n}$.

$H_1: p \leq p_0$ Reject H_0 if $X > c$
 $K_1: p > p_0$ when $P_{p=p_0}(X > c) = \alpha \dots (1)$

We may need to consider a randomized test as binomial distⁿ is discrete and there may not exist an integer c for which (1) will be satisfied.

When n is large, we can consider normal approximation

So here basically the model is that we have x following binomial, say n, p distribution where p is known. Now when n is small, then we can consider test based on x . For example, I can call, say let us define say $P = x/n$. Suppose my hypothesis testing problem is $p \leq p_0$, against $p > p_0$. Then we can consider the test as reject H_0 , if x is $>$ some c where probability of $x > c$ when $p = p_0 = \alpha$.

Now in this case, what will happen is that it is not necessary that we will get exactly $= \alpha$. So we may need to randomize here. We may need to consider a randomized test as binomial

distribution is discrete and there may not exist an integer c , for which one will be satisfied. Now when n is large, we can consider normal approximation.

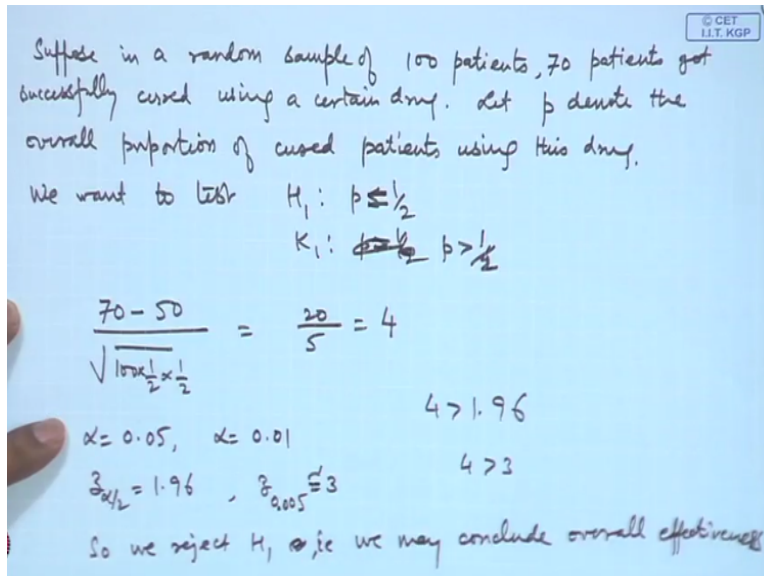
(Refer Slide Time: 03:19)

We can take
$$B_1 = \frac{X - np_0}{\sqrt{np_0q_0}}$$
. When $p = p_0$ and $n \rightarrow \infty$, then
 B_1 converges to $Z \sim N(0, 1)$.
We may consider test based on z_α values
i.e. Reject H_0 when $B_1 > z_\alpha$
Similarly, we may consider $H_2: p \geq p_0$ vs $K_2: p < p_0$
Test is Reject H_2 if $B_1 < -z_\alpha$.
For $H_3: p = p_0$ vs $K_3: p \neq p_0$, then
Test is: Reject H_3 if $|B_1| \geq z_{\alpha/2}$.

We can consider $x - np_0 / \sqrt{np_0q}$. Let us call it say B_1 . Then $p = p_0$ and n tends to infinity, then B_1 converges to z following normal $0, 1$ distribution. Therefore, we can for testing about H_1 versus K_1 , for example this hypothesis, we may consider test based on z alpha values, that is reject H_0 when B_1 is $> z$ alpha. Similarly, we may consider H_2 that p is $\geq p_0$ versus k_2 , $p < p_0$. Then test is reject H_2 if $B_1 < -z$ alpha.

If I consider the hypothesis $p = p_0$ against p is not $= p_0$, then test is reject H_3 if modulus of B_1 is $\geq z$ alpha/2.

(Refer Slide Time: 05:18)



Let me give a simple example. Suppose in a random sample of 100 patients, 70 patients got successfully cured using a certain drug. Let p denote the overall proportion of cured patients using this drug. We want to test say $H_1, p \leq 1/2$ against say $K_1 p > 1/2$ or we may say $p = 1/2$ against $p > 1/2$. Suppose we want to test that the overall effectiveness is more than 50%. In that case, the test statistic will become, you will have $70 - 50 / \text{root } 100 \cdot 1/2 \cdot 1/2$, so that becomes $20/5 = 4$.

So if I consider say $L5 = 0.05$ or $L5 = 0.01$ etc. Then we see that $z_{\alpha/2}$, for example here it is 1.96 and so on. So certainly here $4 > 1.96$. Similarly, at this 1, suppose I say 0.005 then that is still higher value, it is approximately 3, that is $>$ this. We reject H_1 that is we may conclude that overall effectiveness of drug is more than 50%.

(Refer Slide Time: 08:22)

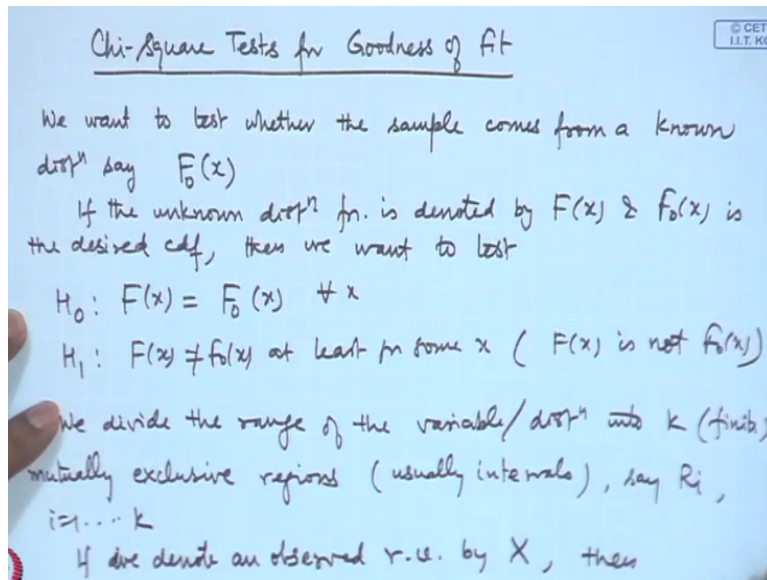
$$\begin{aligned}
 & X \sim \text{Bin}(m, p_1), \quad Y \sim \text{Bin}(n, p_2) \\
 & m, n \text{ are large} \\
 & H_1: p_1 \leq p_2 \quad H_2: p_1 \geq p_2 \quad H_3: p_1 = p_2 \\
 & K_1: p_1 > p_2 \quad K_2: p_1 < p_2 \quad K_3: p_1 \neq p_2 \\
 & \hat{p}_1 = \frac{x}{m}, \quad \hat{p}_2 = \frac{y}{n}, \quad \hat{p} = \frac{x+y}{m+n} \\
 & B_2 = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{m} + \frac{1}{n}\right)}} = \sqrt{\frac{mn}{m+n}} \cdot \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\hat{p}(1-\hat{p})}} \\
 & \text{When } p_1 = p_2, \quad B_2 \text{ has asymptotically } N(0,1) \text{ dist.} \\
 & \text{So we can construct tests for } H_1, H_2, H_3 \text{ based on } B_2
 \end{aligned}$$

Sometimes we may be interested in comparing 2 proportions. That means we have x following binomial m, p_1 and y following binomial n, p_2 and n and m are large. We may need to compare p_1 and p_2 , so we can consider hypothesis of the nature this, or say H_3 . H_2 say $p_1 \geq p_2$ against $K_2, p_1 < p_2$, H_3 $p_1 = p_2$ against K_3 is $p_1 \neq p_2$. So let us refine say $p_1 \hat{=} \text{say } x/m, p_2 \hat{=} \text{say } y/n, p \hat{=} \text{let us define to be } (x+y)/(m+n)$.

And let us define the statistic $(\hat{p}_1 - \hat{p}_2) / \sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{m} + \frac{1}{n}\right)}$ that is actually $= \sqrt{\frac{mn}{m+n}} \cdot \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\hat{p}(1-\hat{p})}}$. So when $p_1 = p_2$, then B_2 has asymptotically normal $0, 1$ distribution. So we can construct tests for H_1, H_2, H_3 , etc. based on B_2 . For example, for H_1 versus K_1 the rejection region will be for $z > z_{\alpha}$. For H_2 versus K_2 , the rejection region will be for $z < -z_{\alpha}$ and for H_3 versus K_3 , the rejection region will be for modulus $z \geq z_{\alpha/2}$, if I am considering Leville α tests.

Let me also consider another related topic. For example, here we are considering in the binomial 2 categories. So for example, if I am considering one binomial, then it is p and then $1-p$ as the proportions of the 2 types. Here we are considering p_1 and p_2 . Now in general we can consider more categories, so this gives actually rise to a test called goodness of fit tests.

(Refer Slide Time: 11:43)



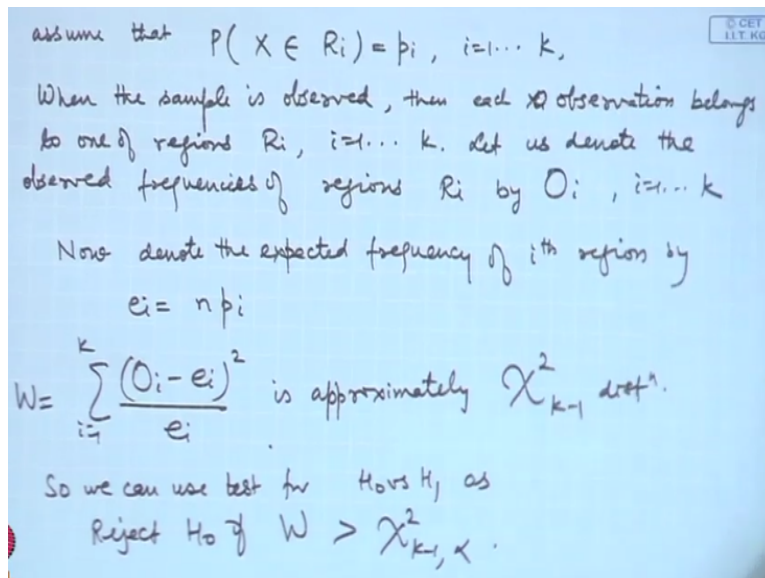
Since asymptotical distributions are Chi square, so the tests are based on that. So we call them Chi square tests for goodness of fit. Let me introduce the problem first. So we want to test whether the sample comes from a known distribution, say $F_0(x)$. In the previous problems, in the usual parametric methods, what we are considering is that we are assuming the form of the distribution, like normal distribution, binomial distribution or I have also given the examples of say, exponential distribution or Poisson distribution.

But there can be situations where we would like to test whether we will have a particular distribution, say binomial distribution or uniform distribution or a Poisson distribution, etc. In that case, we will say that the sample comes from a known distribution, say $F_0(x)$. So you want to test that, that means if the unknown distribution function is denoted by $F(x)$ and $F_0(x)$ is the desired CDF, then we want to test $H_0: F(x) = F_0(x)$ for all x against $H_1: F(x) \neq F_0(x)$ at least for some x .

So that means we are saying that alternative hypothesis is that $F(x)$ is not $F_0(x)$. It could be some other distribution or it may not be distribution. In the Chi square test for goodness of fit, we divide the range of the variable or distribution into k mutually exclusive regions, usually it will be intervals. I mentioned regions, because suppose I am considering binomial distribution, etc., then you have values 0, 1 to n or you are considering Poisson, then it is 0, 1, 2, 3, and so on.

So you will have in finite number of values, but when you can take a practical consideration by considering values by clubbing some of the values together and make it a finite number, so this k is finite. So we divide the range of this, that means we are actually getting some k regions, such as we can give some name here, say R_i $i=1$ to k and if we denote an observed random variable by x then assume that probability of x belonging to the region R_i is some p_i , $i=1$ to k .

(Refer Slide Time: 15:44)



Now what we consider when the sample is observed, then each x_i , each observation belongs to one of regions R_i , $i=1$ to k . Let us denote the observed frequencies of region R_i/O_i for $i=1$ to k . So now we consider suppose n observations are there, we denote the expected frequency of i -th region by $E_i=n \cdot p_i$. So what we do, we construct $\sum (O_i - E_i)^2 / E_i$, $i=1$ to k . This let us call it W , then this has approximately Chi square distribution on $k-1$ degrees of freedom.

So we can use test for H_0 versus H_1 as reject H_0 if W is $>$ Chi square $k-1$ alpha. Let us consider an example here.

(Refer Slide Time: 18:29)

Example: It is assumed that students' preferences for various disciplines are uniformly dist^d. Let there be five options say CS, ECE, EE, ME & CH and let the preference probabilities of these options be p_1, p_2, p_3, p_4, p_5 respectively.

Then we want to test $H_0: p_i = \frac{1}{5}, i=1, \dots, 5$
 $H_1: \text{not so.}$

A random sample of 300 students was taken and their preferences recorded as below

Discipline	CS	ECE	EE	ME	CH	Total
O_i	88	65	52	55	40	300

It is assumed that student's preferences, it is assumed that student's differences for various disciplines are uniformly distributed. So let there be 5 options say CS, EC, EE, ME and CH and let the preference probabilities of these options be say $p_1, p_2, p_3, p_4,$ and p_5 respectively. Then we want to test, that is $p_i = 1/5$ for $i=1-5$ against not so. That means we are assuming the discrete uniform distribution for the preferences.

Then a random sample of say 300 students was taken and their preferences recorded as below. So here we have the branches and the observed frequency O_i is given by 88, 65, 52, 55, and 40. So we want to test whether the preferences are uniformly distributed or not. So we consider here E_i 's. E_i 's are the probabilities of each group. So you notice that, the expected frequency of each group, so if total number is 300, we are assigning probability $1/5$ to each group. So the expected frequency will be 60.

(Refer Slide Time: 22:02)

© CET
I.I.T. KGI

$$\begin{aligned}
 W &= \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} = \frac{\sum (O_i^2 + E_i^2 - 2O_i E_i)}{E_i} \\
 &= \sum \frac{O_i^2}{E_i} + \sum E_i - 2 \sum O_i \\
 &= \left[\sum \frac{O_i^2}{E_i} - N \right] \quad \sum E_i = \sum O_i = N \\
 &= \frac{88^2 + 65^2 + 52^2 + 55^2 + 40^2}{60} - 300 \approx 21.6 \\
 \chi_{4,0.01}^2 &\approx 13.28, \quad \chi_{4,0.05}^2 \approx 9.49 \\
 \text{Hence } H_0 &\text{ is rejected; i.e. students' preferences are} \\
 &\text{biased towards different disciplines.}
 \end{aligned}$$

So we consider here W that = $\sum_{i=1}^k \frac{O_i - E_i}{E_i}$, $i=1$ to k . this is also having an alternative representation. If I expand this numerator, I get O_i square + E_i square - $2O_i E_i / E_i$ that is = $\sum O_i$ square / E_i + $\sum E_i$ - twice $\sum O_i$ = $\sum O_i$ square / E_i + $-N$, because $\sum E_i$ and $\sum O_i$ both equal to the total sample size. So this is an alternative formula for this, so we calculate here by $60 - 300$. So you can do the calculations, it turns out it is = 21.6.

Now there are here 5 groups, so we need to look at Chi square value on 4 degrees of freedom. For example, we may consider say at 0.01 level, then it is 13.28, suppose we consider Chi square value at say 0.5, then it is equal to 9.49. So you can easily see that H_0 is rejected. That student's preferences are biased towards different disciplines.

(Refer Slide Time: 24:27)

If $F_0(x)$ is not completely known, e.g. it may contain unknown parameters $\theta = (\theta_1, \dots, \theta_m)$. In such cases, we have to estimate them from the sample. Consequently, the asymptotic distⁿ of W will be χ^2_{k-m-1} .

Example: 30 randomly selected documents of equal size are taken and the number of typographical errors in them are recorded. The data is summarized below:

No. of Errors	No. of documents
0	6
1	5
2 or 3	8
4 or 5	6
more than 5	5

We want to test whether a Poisson distⁿ appropriately fits the data on no. of errors.

In this particular case, I assume that F_0 completely known. If F_0 is not completely known, for example it may contain, for example I say it is binomial distribution, then there will be an unknown parameter p , which has to be estimated. Suppose, we say it is a Poisson distribution, then the parameter λ has to be estimated. Suppose we say it is normal μ σ^2 distribution, then μ σ^2 have to be estimated first and then they have to be used in the calculation of the expected frequencies.

In that case, the degrees of freedom of the Chi square will be reduced by the number of unknown parameters that have to be estimated from the sample. So it may contain unknown parameters, say $\theta = \theta_1, \theta_2, \dots, \theta_n$. In such cases, we have to estimate from the sample. Consequently, the asymptotic distribution of W will be Chi square $k-n-1$. Let us take one example here.

30 randomly selected documents of = size are taken and the number of typographical errors in them are recorded. The data is summarized below. So if I make a frequency table, number of errors, it is recorded like this 0, 1, then 2 or 3 errors, 4 or 5 errors and more than 5 errors. Then, number of documents who had no errors, it was found to be 6, number of documents which had 1 error were 5, number of documents which had 2 or 3 errors was 8.

Number of documents which had 4 or 5 errors were 6 and the number of documents which had more than 5 errors were 5. We want to test whether a Poisson distribution appropriately fits the data. Because here it is a number of counts, errors or counts. So the data on number of errors, now naturally if we assume, so we have to assume a Poisson lambda distribution. Assume x that is the number of errors follows Poisson lambda distribution. Then this lambda has to be estimated first from the given data.

(Refer Slide Time: 29:01)

The image shows a handwritten derivation on a blue background. It starts with the estimation of lambda as the sample mean: $\bar{x} = \frac{95}{30} = 3.1667$. Then, the Poisson distribution formula is given as $P(X=k) = \frac{e^{-\bar{x}} (\bar{x})^k}{k!}$. Below this, five regions are defined and their probabilities calculated:

- $\hat{p}_1 = \hat{P}(X=0) = \hat{P}(X \in R_1) = e^{-\bar{x}} = 0.04214$
- $\hat{p}_2 = \hat{P}(X=1) = \hat{P}(X \in R_2) = \bar{x} e^{-\bar{x}} = 0.13346$
- $\hat{p}_3 = \hat{P}(X=2) + \hat{P}(X=3) = \hat{P}(X \in R_3) = \frac{\bar{x}^2 e^{-\bar{x}}}{2!} + \frac{\bar{x}^3 e^{-\bar{x}}}{3!} = 0.42434$
- $\hat{p}_4 = \hat{P}(X=4) + \hat{P}(X=5) = \hat{P}(X \in R_4) = 0.28841$
- $\hat{p}_5 = \hat{P}(X > 5) = \hat{P}(X \in R_5) = 0.10164$

 Finally, it states: Then $e_i = np_i = 30 p_i, i=1 \dots 5$. A small logo in the top right corner reads '© CET I.I.T. KGP'.

So we consider this, we will first estimate lambda. So we may consider say maximum likely you would estimate μ or the method of moments estimator. In the case of Poisson distribution, all of them are the same. It is simply \bar{x} . So here you can see it will be equal to simply $95/30=3.1667$. Now based on this, we have distribution written as $e^{-\bar{x}}$ to the power k/k factorial. That is the probability of $x = k$.

So now for example what is probability of $x = 0$. See these are the groups here, like I mentioned here in the very first one that this one that we divide into k mutually exclusive regions here. So k mutually exclusive regions here will correspond to, this is region 1, this is region 2, this is region 3, this is region 4, and this is region 5 here. So what is the probability of region 1, that is probability of region 1. What is the probability that x belongs to region 1.

This is my P1, so that is $= E$ to the power $- \bar{x}$, which of course can be calculated to be 0.04 to 14. Similarly, we can calculate P2 that is probability of $x = 1$, that is the probability of region 2 $= \bar{x} * E$ to the power $- \bar{x}$. One can evaluate it; it turns out to be 0.13346. Now P3 will be the probability of $x = 2$ + probability of $x = 3$, that is the probability of third region, that is $= \bar{x}^2 E$ to the power $- \bar{x}/2$ factorial + $\bar{x}^3 E$ to the power $- \bar{x}/3$ factorial that is $= 0.4343$ etc.

Similarly, probability $x = P4$ that is $= x = 4$ + probability $x = 5$, that is the probability of 4th region, so that will be turning out to be point 28841. That is probability $x = 5$, that is the probability of 5th region that is $= 0.10164$. Now based on this, we can calculate E_i 's are nothing but $n p_i$ that is $= 30 * p_i$, $i=1$ to 5 and we calculate then. So this I can call P1 hat, P2 hat, P3 hat, P4 hat because these are the estimates of the probabilities of these regions. These are the estimates here.

So we calculate W that will be $= \sum O_i^2 / E_i - n = 21.99$. Now Chi square value, you can see, so how many degrees of freedom will be there. We have 5 classes and 1 parameter has been estimated, so it will be 3, so one can easily check the values at some particular level of significance. For example, even at 0.005, it is 12.838, so H_0 is rejected that is the error count do not fit a Poisson distribution. Let me give one more example where Poisson distribution will actually fit the given data.

(Refer Slide Time: 34:19)

Example: The following data represents the frequency count of violent crimes reported in a month for 200 randomly selected districts across a country.

No. of violent crimes	0	1	2	3	4	≥ 5
No. of towns (freq)	22	53	58	39	20	8

We want to test whether the crime count data fits a Poisson distⁿ. $\bar{x} \approx 2$.

e_i | 27 54.2 54.2 36 18 10.6

$W = \sum \frac{O_i^2}{e_i} - n = 2.33$ $\chi_{4, 0.05}^2 = 9.49$

Certainly we have reasons to believe that Poisson

The following data represents the frequency count of violent crimes reported in a month for 200 randomly selected districts across a country. So number of violent crimes and we are clubbing 0, 1, 2, 3, 4, and more than or = 5. So again we would like to test whether it is a Poisson distribution or not. Number of towns, that the frequency, so 22, 53, 58, 39, 20 and 8. So we want to test whether the crime count data fits a Poisson distribution.

So once again, you can check here that \bar{x} is approximately 2, it will be 2. something, so I am just writing 2 here, because that is sufficient for our purpose and we calculate the expected frequencies, expected frequencies will become 27, 54.2, 54.2, 36, 18 and 10.6. So if you calculate W that is $\sum \frac{O_i^2}{e_i} - n$, then that is turning out to be 2.33. So if we look at Chi square value, now since there are 1, 2, 3, 4, 5, 6 groups are there.

The degrees of freedom will be $6-1-1$ and let us take 5% level then it is turning out to be 9.49, so certainly we have reasons to believe that Poisson distribution adequately represents this frequency distribution. Now if we see this thing, the fitting of a distribution problem is basically reducing to a sort of multinomial problem because you are dividing the entire categorized data into k categories.

Now if we are dividing into several categories, then it is immaterial whether we divide it into 1 dimension or we can go for higher dimension also.

(Refer Slide Time: 38:26)

Testing for Independence in $r \times c$ contingency tables

B \ A	A ₁	A ₂	...	A _c	Totals
B ₁	O ₁₁	O ₁₂	...	O _{1c}	O _{1.}
B ₂	O ₂₁	O ₂₂	...	O _{2c}	O _{2.}
...
B _r	O _{r1}	O _{r2}	...	O _{rc}	O _{r.}
Totals	O _{.1}	O _{.2}	...	O _{.c}	N

Observed frequency of $(i,j)^{th}$ cell is denoted by O_{ij}

$$O_{i.} = \sum_{j=1}^c O_{ij}$$

$$O_{.j} = \sum_{i=1}^r O_{ij}$$

Assume theoretical probabilities of $(i,j)^{th}$ cell to be π_{ij}
 Then the marginal probabilities of i^{th} row is $\pi_{i.} = \sum_j \pi_{ij}$
 of j^{th} col is $\pi_{.j} = \sum_i \pi_{ij}$

So let us consider in general testing for independence in r/c contingency tables. So if we are considering contingency tables, then we are considering the classification according to 2 categories A and B and for A, we have categories A₁, A₂, A_r A_c and for B we have B₁, B₂, ...B_c. Now we can actually divide the entire frequency into several cases. Let us put r here. The observed frequencies, I am writing as O₁₁, O₁₂, O_{1c}, O₂₁, O₂₂, O_{2c} and O_{r1}, O_{r2}, O_{rc}.

We consider the row and column sums. So if we sum the first row, we call the sum as O_{1 dot}, O_{2 dot} and so on or dot. Similarly, if we sum the columns, we call that O dot 1, O dot 2, and so on, O dot c. The total sum is n. So we have the following notations observed frequency of ij-th cell is denoted O_{ij} and then we define O_{i dot} that is = sigma O_{ij} for j=1 to c, so simply the summations and similarly O dot j that = sigma O_{ij}, i=1 to r. These are the row and column totals.

Then if we are assuming that the 2 things are independent, there will be theoretical probability of assume theoretical probabilities of ij-th cell to be pi ij, then the marginal probabilities of i-th row is pi I dot that = sigma pi ij sum over j and of j-th column, it is pi dot j that = sigma pi ij sum over i.

(Refer Slide Time: 42:26)

If the row & column are independent then we must have

$$\pi_{ij} = \pi_{i.} \times \pi_{.j}$$

So we calculate the expected freq of $(i,j)^{th}$ cell using this assumption

$$e_{ij} = \frac{O_{i.} \times O_{.j}}{N}, \quad N = \sum \sum O_{ij}$$

$$W^* = \sum_i \sum_j \frac{(O_{ij} - e_{ij})^2}{e_{ij}}$$

has asymptotically $\chi^2_{(r-1)(c-1)}$ dist.

We will reject the hypothesis of independence if

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

If the row and columns are independent, then we must have $\pi_{ij} = \pi_{i.} \times \pi_{.j}$. So we calculate the expected frequency of ij -th cell using this assumption. So that is $e_{ij} = O_{i.} \times O_{.j} / N$. Here N is actually the sum of all the frequencies. So if I use this, then we get, let us call it $W^* = \text{double summation } (O_{ij} - E_{ij})^2 / E_{ij}$. This has asymptotically Chi square $r-1$ $c-1$ distribution. So we will reject the hypothesis of independence, if W^* is $>$ Chi square $r-1$ $c-1$ alpha.

(Refer Slide Time: 44:40)

Example: The following data represents the number of accidents taking place in 3 shifts of 4 factories producing an item. The data is recorded for a year. We want to test whether the incidence of accidents is independent type of factory & shift

Shift \ Factory	A	B	C	D	Totals
Shift 1	10	12	6	7	35
Shift 2	10	24	9	10	53
Shift 3	13	20	7	10	50
Totals	33	56	22	27	$N=138$

$W^* \approx 1.81$

$\chi^2_{6, 0.05} = 12.59$

We can say that shift & factories are independent as $W^* <$ $\chi^2_{6, 0.05}$

$e_{11} = \frac{33 \times 35}{138}$ $e_{23} = \frac{22 \times 53}{138}$,

Let me give one application here. The following data represents the number of accidents taking place in 3 shifts of 4 factories producing an item. The data is recorded for a year. So we want to test whether the incidence of accidents is independent, that means whether in a particular factory

at particular shift has more accidents are less, so independent of type of factories and shifts. So the data is recorded in this particular fashion.

Suppose we have 4 factories A, B, C, D and the data is recorded over shift 1, shift 2 and shift 3. That is 10, 10, 13, 12, 24, 20, 6, 9, 7, 7, 10, 10. If we consider the totals, this is 33, this is 56, 22, and 27 and on this side, if we consider the row totals, it is 35, 53, 50, the total $N = 138$. So we calculate for example, what will be E_{11} . E_{11} will be $33 \cdot 35 / 138$. Similarly suppose, I consider say E_{23} , so E_{23} will be $22 \cdot 53 / 138$ etc.

So we calculate the W^* that is turning out to be here 1.81 approximately. Now if I consider Chi square on $2 \cdot 3$ that is 6 degrees of freedom at a particular level, say 0.05, then it is 12.59. So we can say that shifts and factories are independent with respect to occurrence of accidents. You can say that the incidence of accidents is homogenous across the factories. Let us take 1 or 2 more applications of the testing and these problems.

(Refer Slide Time: 49:00)

Example: Over two seasons a professional player (basketball player) was at field exactly 5 minutes in about 200 games.
 $X_i \rightarrow$ the no. of hits he makes in game i , $i = 1, \dots, 200$.
 $X_i \rightarrow 0, 1, 2, 3, 4$ (assume)

Value of X_i :	0	1	2	3	4
no. of X_i 's:	73	82	38	7	0

We want to test whether a binomial distⁿ. will fit the data

$$p_1 = P(X=0) = (1-p)^4, \quad p_2 = P(X=1) = 4p(1-p)^3, \quad p_3 = P(X=2) = 6p^2(1-p)^2$$

$$p_4 = P(X=3) = 4p^3(1-p), \quad p_5 = P(X=4) = p^4$$

$$L(p) = \frac{200!}{73! 82! 38! 7! 0!} (1-p)^{73} (4p(1-p))^82 (p^4)^0 (6p^2(1-p))^38 (4p^3(1-p))^7$$

$L(p)$ is maximized when $\hat{p} = 0.224$.

Over 2 seasons, a professional player of some game, we may consider for example a basketball player exactly 5 minutes in about 200 games. So x_i is the number of hits he makes in game i , $i=1$ to 200. Each x_i can take value 0, 1, 2, 3, 4. So we have the following data, value of x_i is 0, 1, 2, 3, 4 and number of x_i is 73, 82, 38, 7, 0. We want to test whether a binomial distribution will fit the data. Now in a binomial distribution, we have a parameter p here.

So let us consider say p hat. Based on this data, we can calculate actually. So P_1 that is probability $x = 0$ that is $= 1-p$ to the power 4, P_2 that is probability $x = 1$ that $= 4p * 1-p$ cube, P_3 that $=$ probability $x = 2$ that $= 6p$ square $* 1-p$ square, P_4 that is probability $x = 3$ that is $= 4p$ cube $* 1-p$ and $P_5 =$ probability $x = 4$ that $= p$ to the power 4.

So we have the likelihood function that is 200 factorial / 73 factorial, 82 factorial, 38 factorial, 7 factorial, 0 factorial $* 1-p$ to the power 4 to the power 73 $4p * 1-p$ cube to the power 82 p to the power 4 to the power $0 * 6 p$ square $* 1-p$ square to the power $38 * 4 p$ cube $* 1-p$ to the power 7 . So this can be simplified L hat p is L_p is maximized when $p=0.224$. So based on this, we can calculate P_1 hat that is 0.363 , P_2 hat $= 0.419$, P_3 hat $= 0.181$, P_4 hat $= 0.035$, P_5 hat $= 0.003$ etc.

So if you calculate this, calculate the Chi square value here that $= 0.178$ approximately. So if you compare with Chi square value on here we have 5 categories and 1 parameter has been estimated, so you will have it on 4 degrees of freedom and one can see this. I will give one application of the general testing problem, which we discussed for the normal populations.

(Refer Slide Time: 54:11)

Testing Example for Normal Populations

$n_1 = 121, \bar{x}_1 = 2.6, s_1^2 = 1.44$
 $n_2 = 61, \bar{x}_2 = 0.4, s_2^2 = 0.0121$

To test equality of means, we need to firstly test the equality of variances

$H_0: \sigma_1^2 = \sigma_2^2$
 $H_1: \sigma_1^2 \neq \sigma_2^2$

$F = \frac{S_1^2}{S_2^2} \approx 119.0$
 $F_{121, 60, 0.1} = 1.34$

So H_0 is rejected.

$H_0^*: \mu_1 = \mu_2$
 $H_1^*: \mu_1 > \mu_2$

$T = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \approx 20.0$

$v = 123$, Certainly H_0^* is rejected.

Testing example for normal populations. The summary data is given by, we have 2 samples for 2 types of elements present in the bones of children and then the following data is collected, n_1 is 121, $\bar{x}_1 = 2.6$, $s_1^2 = 1.44$, $n_2 = 61$, $\bar{x}_2 = 0.4$, $s_2^2 = 0.0121$. We want to test

whether the 2 normal populations have similar means or variances. See if you calculate this, firstly we test to test equality of means. We need to firstly test the equality of variances.

So that means, we test say $H_0 \sigma_1^2 = \sigma_2^2$ against $H_1, \sigma_1^2 \neq \sigma_2^2$. Let us calculate the statistic s_1^2 / s_2^2 and it turns out to be 119.00 approximately. So if I consider say F on 120, 60 degrees of freedom, then the values, say at 0.1 that will 1.34 etc. This is certainly larger. So H_0 is rejected. So now I consider say $\mu_1 = \mu_2$ against say $\mu_1 > \mu_2$.

Then we formulate the test statistic $(\bar{x}_1 - \bar{x}_2) / \sqrt{s_1^2/N_1 + s_2^2/N_2}$. That is approximately 20.0. If I consider the degrees of freedom of this T here that is turning out to be approximately 123, so certainly H_0 is rejected. So here for testing the equality of the means, which procedure is to be used, because I discussed 4 different procedures, firstly we need to check about variance.

Now for the variance here it turns out that it is rejected here and therefore we have followed this procedure. If it was accepted, then we have to follow another one, which was based on the pooling procedure. So depending upon what actual method will be used, then only you apply the testing methodology. We have discussed some of the important parametric methods. There are many more, but in this particular course, I will restrict attention to this.

In the following lectures, I will move over to multivariate analysis. So we will have elementary discussion of the multivariate normal distribution and then the related distributions and how they are used for certain calculations or computations or inferences when you have multivariate data. So in that following lectures, we will take up that.