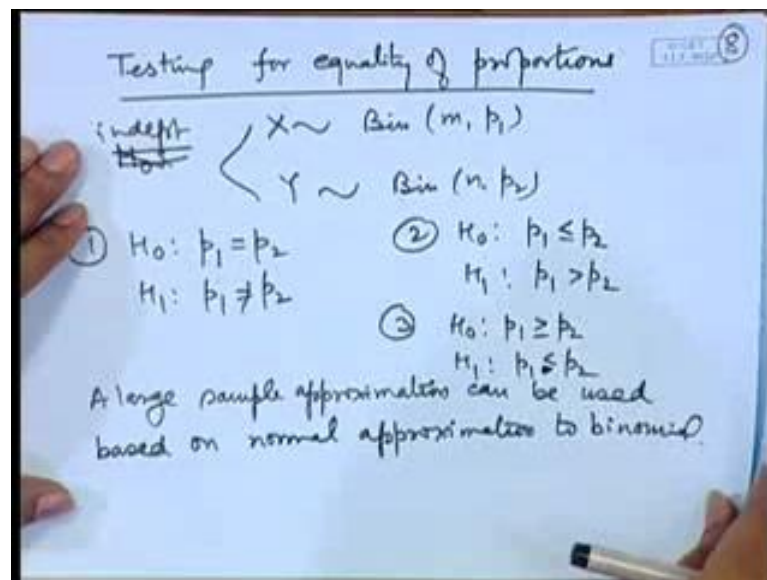


Probability and Statistics
Prof. Somesh Kumar
Department of Mathematics
Indian Institute of Technology, Kharagpur

Lecture – 75
Testing Equality of Proportions

We also have the normal test when we are comparing the proportions of the 2 binomial populations. So, for example, the data is recorded not in the numerical measure, but in the characteristic form. For example, we want to see certain opinions; we want to see the effect of certain say learning procedure. So, that may be the result may be in the form of that a test is conducted, so for example, a set of a student is start a certain instructional material, another set of a student start another instructional material, a common test is conducted we want to see how many passed in the first set and how many passed in the second. Is there a significant difference in the proportions? That means, we want to see whether the instructional material one is better or the instructional material two is better.

(Refer Slide Time: 01:23)



So, basically the situation is of the following form testing for equality of proportions. So, the general model is of the say X follows say binomial m, p 1 and Y follows binomial n, p 2.

We assume these 2 samples to be independently taken here. So, we are interesting in testing hypothesis of the form p 1 is equal to p 2, against say p 1 is not equal to p 2 say

this is 1. We may test p_1 is less than or equal to p_2 , against p_1 is greater than p_2 or say H_0 : $p_1 \geq p_2$; p_1 is greater than or equal to p_2 , against H_1 : $p_1 < p_2$ etcetera. A large sample approximation can be used based on normal approximation to binomial.

(Refer Slide Time: 03:03)

$$\hat{p} = \frac{x+y}{m+n}, \hat{p}_1 = \frac{x}{m}, \hat{p}_2 = \frac{y}{n}$$

$$Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{m} + \frac{1}{n}\right)}} = \sqrt{\frac{mn}{m+n}} \cdot \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})}}$$

Under $p_1 = p_2$, $Z \sim N(0,1)$ for m & n large.

For ① Reject H_0 if $|Z| \geq z_{\alpha/2}$
 For ② Reject H_0 if $Z \geq z_{\alpha}$
 For ③ Reject H_0 if $Z \leq -z_{\alpha}$.

So, we may calculate here \hat{p} as $\hat{X} + \hat{Y}$ by $m + n$, \hat{p}_1 is equal to X by m , \hat{p}_2 is equal to Y by n . So, this is the first proportion, this is the second proportion sample proportion, this is the pooled proportion. So, we construct a statistic a Z as \hat{p}_1 minus \hat{p}_2 divided by square root \hat{p} into $1 - \hat{p}$, 1 by $m + 1$ by n , which is actually equal to root of $m n$ by $m + n$, \hat{p}_1 minus \hat{p}_2 divided by root \hat{p} into $1 - \hat{p}$.

Under the assumption that p_1 is equal to p_2 Z is approximately normal $0, 1$ for m and n large. So, we may make a test based on this for 1 reject H_0 if modulus Z is greater than or equal to $z_{\alpha/2}$. For 2 it will be reject H_0 if Z is greater than or equal to z_{α} . For 3 reject H_0 if Z is less than or equal to minus z_{α} . So, this is an approximate test, if the assumption that m and n are large is not true in that case we may have to go for a exact procedure, but that procedure will make use of the distribution which is calculated from the binomial.

So, under p_1 equal to p_2 , the distribution of $X + Y$ is again binomial and one can make use of the distribution of X given $X + Y$, which is hyper geometric and there is a

test procedure for that, but we are not going to discuss that in this particular course here. Let me give an application of this here.

(Refer Slide Time: 05:27)

Example. Suppose one wants to compare the effectiveness of treatments by two different surgical procedures for a certain disease.

	Treatment 1	Treatment 2
Successfully	63	107
Failure	37	43
	100	150

$p_1 = 0.63$, $\hat{p}_2 = 0.71$, $\hat{p} = 0.68$
 $Z = -1.33$, $Z_{0.05} = 1.645$
 $Z_{0.01} = 2.32$
 $H_0: p_1 = p_2$, H_0 cannot be rejected.
 $H_1: p_1 \neq p_2$

Suppose one wants to compare the effectiveness of treatments by two different surgical procedures for a certain disease.

So, for this set of patients is taken on one set of patients, one surgical procedure is adopted and we observe the proportion of success. So, let us make the data in this particular fashion, suppose 100 patients are there on which treatment procedure 1 is adopted, we see how many are successfully treated and how many of them are failures. Suppose in out of 100 year, it turns out that 63 are successfully treated 37 are unsuccessful whereas, using the treatment procedure 2; suppose 100 patients 150 patients were given this procedure out of that 107 where successfully treated and 43 are not successfully treated; that means, on them the surgical procedure did not yield any positive result.

Let us look at the proportions here p_1 hat is 0.63, p_2 hat is equal to 0.71 and p hat is equal to 0.68. So, if we calculate the Z statistic here, that is p_1 hat minus p_2 hat divided by root of p hat into 1 minus p hat into $m n$ by $m + n$, this value turns out to be minus 1.33. If we are considering any reasonable level of significance say Z 0.05 that is 1.645, suppose I take z 0.01, that is 2.32 etcetera; then you can see that we if we consider say hypothesis H_0 : p_1 is equal to p_2 against say H_1 , p_1 is not equal to p_2 .

Then at level of significance say 10 percent, 2 percent etcetera. H_0 cannot be rejected, because the absolute value of z that is 1.33 is a smaller than this values.

(Refer Slide Time: 09:01)

	Treatment 1	Treatment 2
Success	3	107
Failure		43
Total		150

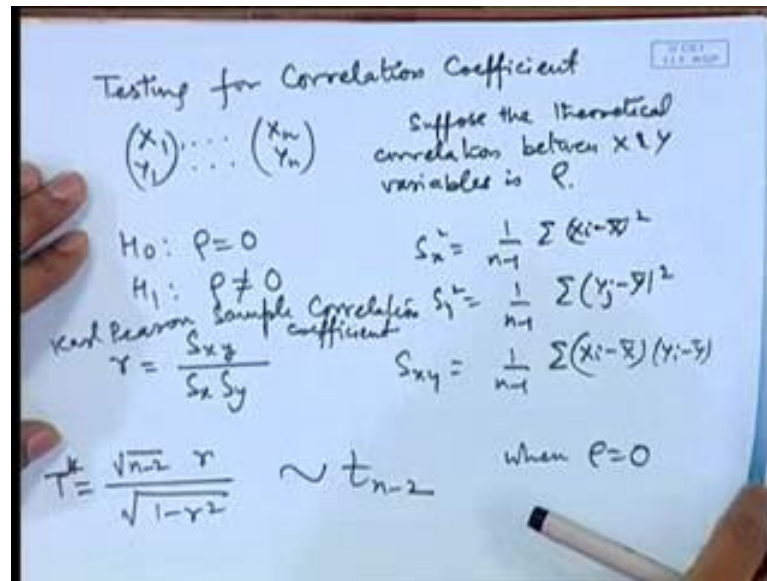
0.71 , $\hat{p} = 0.68$
 $z = 1.33$
 $z_{0.05} = 1.645$
 $z_{0.01} = 2.32$
 $H_0: p_1 = p_2$ H_0 cannot be rejected.
 $H_1: p_1 \neq p_2$ So no significant difference in the success rates of the procedures.

So, that means, there is no significant difference in the 2; in the success rate of the 2 surgical procedures. So, no significant difference in the success rates of 2 procedures. Although from here it looks that the second procedure is more effective, but statistically speaking there is no significance difference here.

We have seen here that sometimes the assumption of the correlated observations is required as in the case of exercise program etcetera, but we have seen the use of t whether t test should be done under certain checks; that means, we should use it with k for example, we are adopting a paired t test procedure, but the value of the correlation is say 0, in that case $\sigma_1^2 + \sigma_2^2$ will be the variance of $X_i - Y_i$; that means, we have unnecessarily used reduced over degree of freedom here, we are using only $n - 1$ degrees of freedom, consequently our power of the test will reduce.

So, it is not advisable to go for a paired t test here; that means, in such a case it will be a reasonable option. Firstly, to check whether the correlation is 0 or not in the given data set. So, fortunately we can find actually a test for the correlation coefficient being significantly different from 0 or not.

(Refer Slide Time: 10:53)



So, testing for correlation coefficient; so we have the data $X_1, X_2, X_n; Y_1, Y_2, Y_n$. So, in the 2 data sets we consider the correlate sample correlation coefficient and the population correlation coefficient suppose the theoretical correlation between X and Y variables is rho.

So, we want to test whether rho is equal to 0, against say rho is not equal to 0. For this we calculate the sample correlation coefficient that is r, that is equal to $S_{xy} / (S_x S_y)$; where this S_x^2 is $1 / (n - 1) \sum (X_i - \bar{X})^2$, S_y^2 is $1 / (n - 1) \sum (Y_j - \bar{Y})^2$ and S_{xy} is equal to $1 / (n - 1) \sum (X_i - \bar{X})(Y_i - \bar{Y})$, that is the sample standard sample variances for the 2 samples and this is the sample covariance.

So, based on this we define the Karl Pearson sample correlation coefficient, then it has been observed that square root $n - 2$ into r divided by the root $1 - r^2$ this is having a t distribution on $n - 2$ degrees of freedom, when rho is equal to 0, let me call it t star. So, one can make use of this for testing the significance of correlation; that means, whether there is a significant correlation between the 2 variables or not.

(Refer Slide Time: 13:26)

Example: Suppose the data on incomes of parents and children is available. It is assumed that there is no significant correlation between the average incomes of parents & their children. 43 pairs are selected.

$$r = 0.412$$
$$T^* = \frac{\sqrt{41} (0.412)}{\sqrt{1 - 0.412^2}} = 2.9.$$

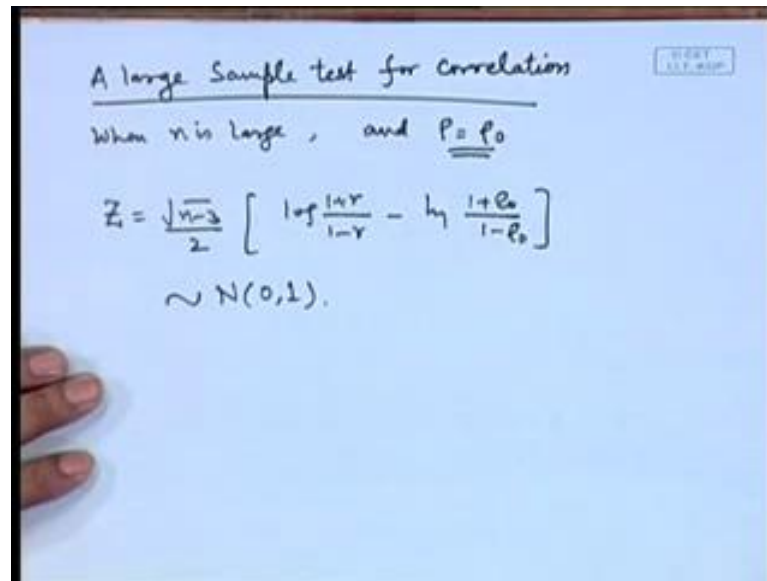
$t_{0.005, 41} = 2.7$. So H_0 is rejected at 1% level. i.e. incomes of children are related to the parents.

Let us consider one example here; suppose the data on incomes of parents and children their children is available, it is assumed that there is no significant correlation between the average incomes of parents and their children. So, 43 pairs are selected and their sample correlation was calculated which turned out to be 0.412. So, from here this T star value that is square root 41 into 0.412 divided by square root 1 minus 0.412 square is calculated, this value turns out to be 2.9.

Suppose we are considering say t on say 0.005 at 41, this value is 2.7. Now you can see here that this is extremely small level of significance we are taking. So, H_0 is rejected at any reasonable level of significance, this is actually 1 percent level. So, at 1 percent level itself this is rejected; that means, the incomes of children are related to the parents; that means, higher income parents their children will also tend to earn higher incomes, and lower income parents their children will have lower incomes.

Now, this test is again based on the normality assumption; sometimes the normality assumption may not be valid, in that case there is a large sample test for correlation coefficient.

(Refer Slide Time: 16:17)



A large Sample test for Correlations

When n is large, and $\rho = \rho_0$

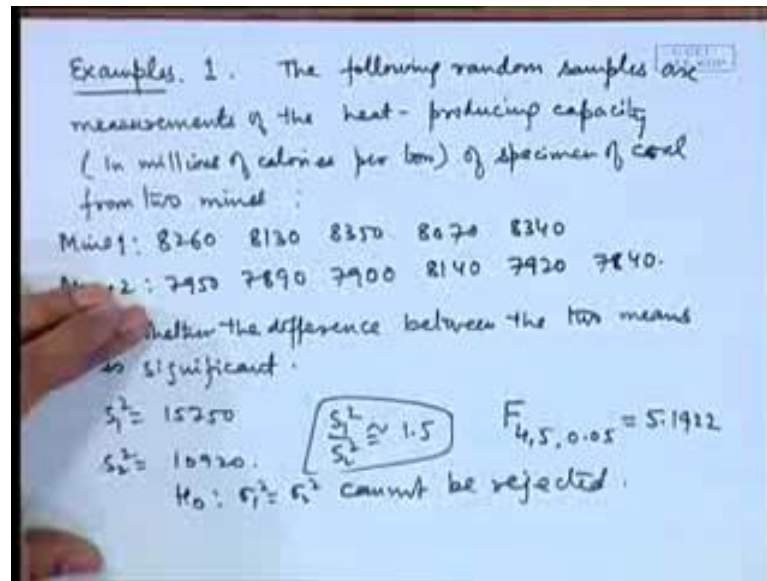
$$Z = \frac{\sqrt{n-3}}{2} \left[\log \frac{1+r}{1-r} - \log \frac{1+\rho_0}{1-\rho_0} \right]$$

$\sim N(0,1)$.

When n is large and say ρ is equal to ρ_0 . So, we may not test that ρ is equal to 0 but some arbitrary value ρ_0 , then we construct Z that is equal to root n minus 3 by 2, log of $1 + r$ divided by $1 - r$ minus log, $1 + \rho_0$ divided by $1 - \rho_0$ then this has approximately normal 0, 1 (Refer Time: 17:17).

So, consequently if n is large, then we may test about the correlation being equal to any arbitrary value. Of course, if ρ is 0 here then this term will vanish and we will have only this particular term. So, this is an approximate normal test here and this is not based on the assumption of normality for the initial samples that is X_i 's and Y_i 's here. Let us take a few more examples here of the test that we have discuss today. So, let us take one example based on say normal distributions here.

(Refer Slide Time: 18:23)



The following random samples are measurements of the heat producing capacity in millions of calories per ton of a specimen of coal from two mines. So, in mine 1: there will be 5 measurements are taken, the values are 8260, 8130, 8350, 8070, 8040. In the mine 2: 6 observations were taken, 7950, 7890, 7900, 8140, 7920 and 7840. Test whether the difference between the 2 means is significant.

So, here 2 samples are available to us, and we are having 5 and 6 observations respectively from the 2 samples. So, what we do here, we have to check whether μ_1 is equal to μ_2 or μ_1 is not equal to μ_2 , but once again here to test this hypothesis we will have to test about the variances also. So, whether the variances are same? So, here we see that if we consider S_1^2 here, S_1^2 is equal to 15750 etcetera S_2^2 is equal to 10920.

So, if we consider the ratio here, S_1^2 by S_2^2 that is approximately 1.5. So, if I am looking at the F value on say 1, 2, 3, 4, 5. So, 4 and 5 degrees of freedom, let us see one example here from the tables of the F distribution and say at 0.05 level, if we are seeing $F_{4,5}$. So, at 0.05 the value is 5.1922. So, $H_0: \sigma_1^2 = \sigma_2^2$ cannot be rejected. Now if this cannot be rejected then for the equality of means, we will go for the pooled sample variance procedure.

(Refer Slide Time: 22:14)

So for testing equality of means we will use pooled sample variance procedure

$$S_p^2 = \frac{4S_1^2 + 5S_2^2}{9} = \frac{4 \times 15750 + 5 \times 10920}{9}$$

$\bar{X} = 8230$
 $\bar{Y} = 7940$

$$S_p = 114.31$$
$$T = \sqrt{\frac{m \cdot n}{m+n}} \frac{\bar{X} - \bar{Y}}{S_p} = 4.19$$

$t_{9, 0.005} = 3.25$

So at 1% level also the hypothesis of equality of means is rejected.

So, for testing equality of means we will use pooled sample variance procedure. So, if we see S_p^2 that is equal to $4S_1^2 + 5S_2^2$ by 9 that turns out to be. So, that is equal to well that is a huge value here, $4 \times 15750 + 5 \times 10920$ by 9. So, that is equal to some value. So, S_p is taken to be square root of that that is 114.31.

So, if we calculate the T variable here that is $\sqrt{\frac{m \cdot n}{m+n}} \frac{\bar{X} - \bar{Y}}{S_p}$. Then first thing is we observe here that \bar{X} is equal to 8230 and \bar{Y} is equal to 7940. So, this value turns out to be after calculation 4.19. So, if we see the t values at say 9 degrees of freedom then even at 0.005 this value is 3.25. So, at 1 percent level also, the hypothesis of equality of means is rejected.

So, we conclude that in the 2 mines, the measurements of the heat producing capacity are significantly different. Because it may be due to difference in the type of the coal that is available, it may be due to the type of the mine that you are having may be in, one of the mines you have a very low level roots and various kind of parameters which may be operating in those mines there. So, it may be due to that. So, to sum up if we are comparing the means of 2 normal populations, the first thing is we have to look at is that what type of variances are there. If the variances are known, then we have one type of procedure. If the variances are unknown, then we firstly test whether the variances are same or not; if they are same then we go for a pooled sample variance procedure, if they are not same then we go for a different procedure which is an approximate test.

We also see the correlation if the correlation is present then we may go for a pairing, if we do not have the correlation then we may go for independent samples. So before adopting any procedure, one has to carefully examine the problem and then choose the appropriate test; we have also seen the effect of choosing the null and alternative hypothesis. As I have already mentioned since we are controlling the probability of type one error therefore, it is always reasonable to put a stronger hypothesis or you can say the convection in which we have more that hypothesis as an alternative hypothesis, because rejection of the hypothesis strong conclusion whereas, acceptance of the hypothesis becomes a weaker conclusion. Simply because of the reason that we are actually concluding the probability of type one error.

In the forth coming lecture I will be discussing the chi square test for goodness of fit or testing for the independence.