

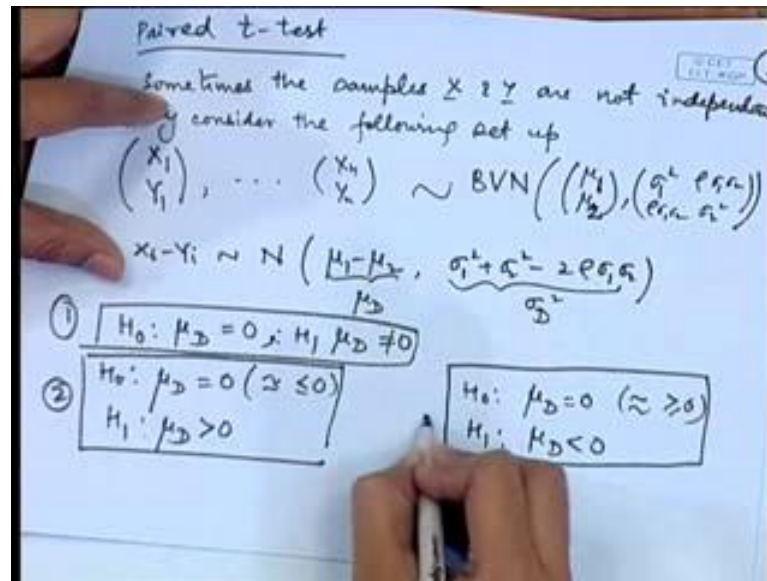
Probability and Statistics
Prof. Somesh Kumar
Department of Mathematics
Indian Institute of Technology, Kharagpur

Lecture – 73
Paired t-Test

There is yet another situation which arises here for example, here we have made an assumption that the 2 samples are taken independently and we are comparing the means here, but there may be situations where the samples may not be independent for example, we are testing the effect of a certain medicine on some patients. So, earlier say one set of patients is chosen, they are given one medicine. Now we observe the effects of that medicine for example it is a medicine to reduce the blood pressures. So, the adverse effectiveness of this medicine is recorded, now the same set of patients is given another medicine may be after one month and then again the effect of the medicine is recorded.

Now, here since the set of the patients is a same therefore, the observations X_i s and Y_i s suppose I call the first set as X_1, X_2, X_n and the second set of observations as Y_1, Y_2, Y_n here we cannot assume them to be independent. Now if we are not able to assume that they are independent then the procedures that I described earlier in case 1, case 2 and case 3 cannot be adopted here, because these procedures assume that the things are independent because we have made use of the additive properties of the chi square distribution or the we have use say linearity property of the normal distribution etcetera.

(Refer Slide Time: 02:05)



So, in this case a simplified procedure is proposed and we call it paired t-test. So, sometimes the samples X and Y are not independent. So, we may consider the following set up. So, we may now consider that X_1, Y_1 these are observations on the same entity. So, $X_1, Y_1, X_2, Y_2, \dots, X_n, Y_n$ this follows a Bivariate normal distribution with mean say μ_1, μ_2 and variance say σ_1^2, σ_2^2 and a certain correlation coefficient say ρ here.

So, in that case and we have to test about the equality of μ_1 and μ_2 or μ_1 less than or equal to μ_2 etcetera. Now once again you see if we have a bivariate model then if we consider the differences that is $X_i - Y_i$'s they again follow univariate normal distribution with means $\mu_1 - \mu_2$ and variance will become $\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2$. We can call it say; so we give a notation here D_i and this I call say μ_D and this is σ_D^2 . So, the testing problem can be reduced to about μ_D , whether μ_D is equal to 0, or μ_D is not equal to 0 or say $H_0: \mu_D = 0$ against say μ_D is greater than 0 of course, this is equivalent less than or equal to 0 here or μ_D is equal to 0 against μ_D less than 0 of course, this is again equivalent to greater than or equal to 0 itself. Because these hypothesis is now equivalent for example, μ_1 is equal to μ_2 , this is equivalent μ_1 not equal to μ_2 , this is equivalent μ_1 less than or equal to μ_2 , this is equivalent to μ_1 greater than μ_2 etcetera.

Now, this model again reduces to the model which we considered for one population; that means, testing for the mean of a normal distribution and we know that there is a t test for that.

(Refer Slide Time: 04:51)

$$T^* = \frac{\sqrt{n} \bar{D}}{S_D}$$
 where $\bar{D} = \frac{1}{n} \sum D_i$, $S_D^2 = \frac{1}{n-1} \frac{\sum (D_i - \bar{D})^2}{\sum (D_i - \bar{D})^2}$

Under $\mu_D = 0$,
 $T^* \sim t_{n-1}$.

Therefore we can have a test based on T^* .

For problem ①, the test will be
 Reject H_0 if $|T^*| \geq t_{n-1, \alpha/2}$

For ②, Reject H_0 if $T^* \geq t_{n-1, \alpha}$

For ③, Reject H_0 if $T^* \leq -t_{n-1, \alpha}$

So, if we consider here the statistic defined by say T star is equal to root n D bar. So, let me write it root n D bar divided by S D; where D bar is the mean of D i's, S D square is the sample variance of the D i minus D bar square, then one can see that under mu D is equal to 0 the distribution of T star will be t on n minus 1 degrees of freedom therefore, we can have a test based on T star. So, for example, if I am considering let me name the hypothesis here let me call this 1 as say hypothesis problem 1, this 1 I can consider as problem 2, this 1 I can consider as say problem 3 ok.

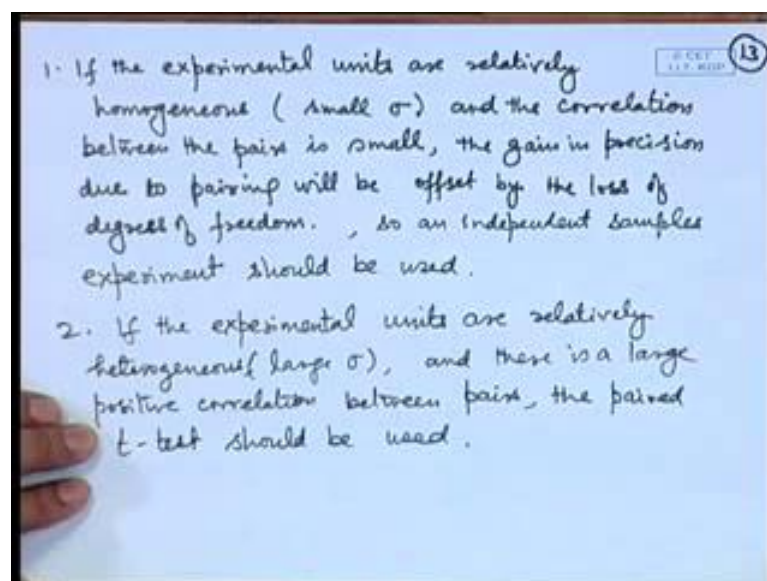
So, for hypothesis testing problem 1, the test will be reject H naught if modulus of T star is greater than or equal to t n minus 1, alpha by 2. For problem 2 it will be reject H naught if T star is greater than or equal to t n minus 1 alpha. For the testing problem 3 the test will be reject H naught, if T star is less than or equal to minus t n minus 1 alpha. Once again let me mention here that this equality or inequality in the less than or equal to or greater than or equal to it does not make any difference, because we are dealing with the continuous random variables, so the probability of the equality is 0. Now the case that I have discussed here that we may have the dependence on the observations, but another problem that we have seen here that there is a loss of degrees of freedom here.

So, if we are losing, degrees of freedom here then the power of the test is also reduced here. So, then one has to make a very clear cut choice that in a given situation whether the data is independent or not. If the data is actually showing dependency; that means, there is a high degree of correlation, then naturally the divisor becomes.

For example if you are considering σ_d^2 . So, if ρ is higher than it is going to affect, because if they are independent then the variances will become $\sigma_1^2 + \sigma_2^2$ for $X_i + Y_i$ or $X_i - Y_i$ if we are taking them independent whereas, here you are making it is smaller.

So, consequently the estimate of that will also become a smaller; that means the test is becoming more sensitive, because if this value is becoming a smaller then when you are taking the ratio in the ratio you are dividing by smaller value that is in this statistic. So, this value is becoming larger; that means, there are more chances of rejection. If we are having more chances of rejection, then the probability of that is $1 - \text{probability of type 2 error}$ that is the power of the test will increase, so the test will become slightly more powerful. On the other hand you are losing certain degrees of freedom here. So, if the correlation is not much and still you are taking the dependence model, then you may have a loss here. So, some general guidelines can be given here let me just write it here.

(Refer Slide Time: 09:35)

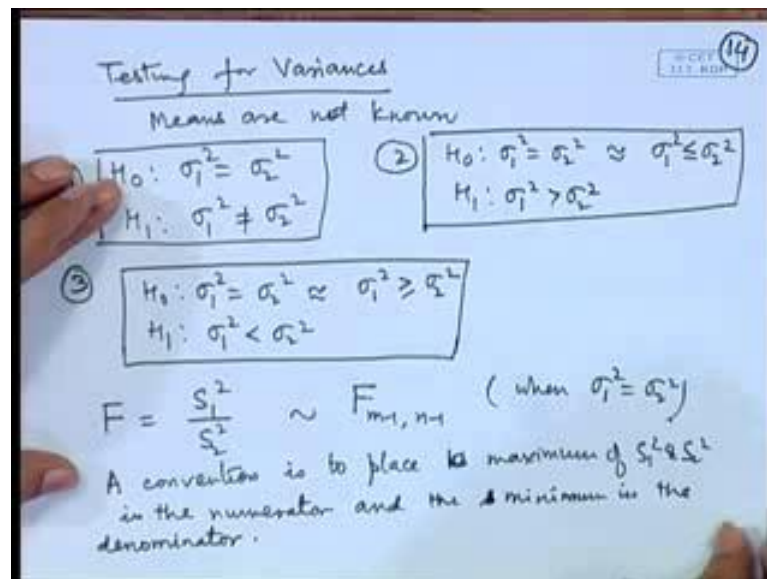


If the experimental units are relatively homogeneous; that means, variability is small and the correlation between the pairs is small, the gain in precision due to pairing will be

offset by the loss of degrees of freedom. So, an independent samples experiment should be used. On the other hand if the experimental units are relatively heterogeneous, that is large sigma and there is a large positive correlation between pairs, the paired t-test should be used.

So, these are some general guidelines, but of course, when a given problem is there one has to see carefully, whether the pairing is permissible or not; if it is permissible then generally there is a better chance that you will have a better outcome or you can say a better result if you use the pairing or you can say paired t-test. Now another important point is that when we were testing for the equality of the means, we considered different cases and they were related to certain information about the variances; that means, naturally there is a question of checking the equality of the variances, because if the variances are unknown, but equal you have another procedure if you have them to be unequal then you have another procedure; that means, one should carry out a test for the equality of the variances, before testing by equality of the means.

(Refer Slide Time: 13:01)



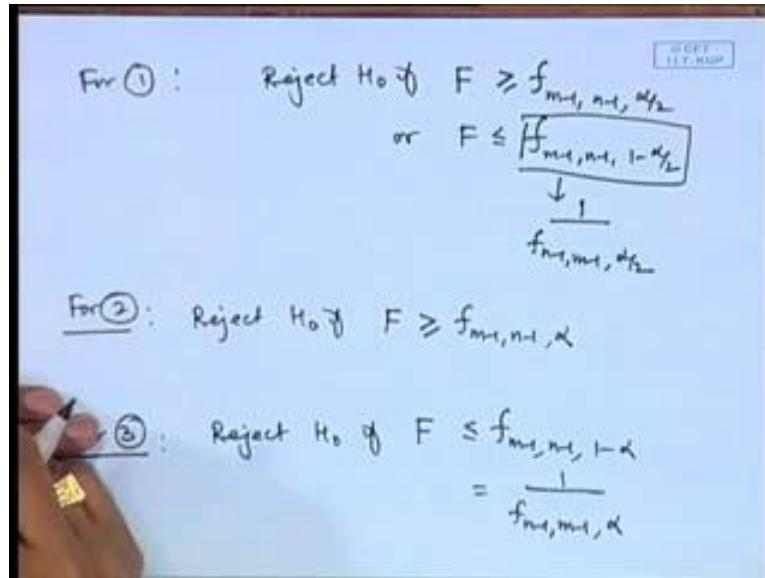
Now, we can produce a test for variances also, testing for the variances of course, once again we may have the case where the means are known or unknown, but it will not make too much difference here. In fact, we have seen in the case of testing for the normal variance that you have a change in the single degrees of freedom there. So, there is not too much sensitivity involved there, but in general as we know that the means may not be

known here, since we are going to test for the equality of the means also. So, let us consider the case means are not known. So, we have σ_1^2 is equal to σ_2^2 against say $H_1: \sigma_1^2 \neq \sigma_2^2$. Once again let me call this as hypothesis testing problem one, another problem would be where you have σ_1^2 is equal to σ_2^2 against $\sigma_1^2 > \sigma_2^2$ and once again this H_0 is equivalent to $\sigma_1^2 \leq \sigma_2^2$.

Another case could be $H_0: \sigma_1^2 = \sigma_2^2$, which is of course, equivalent to $\sigma_1^2 \geq \sigma_2^2$, against $\sigma_1^2 < \sigma_2^2$. So, that is give a test of hypothesis for all the 3 cases. We can consider a F test here. So, F is S_1^2 by S_2^2 of course, the convention is that the bigger value among S_1^2 and S_2^2 should be put in the numerator, the reason is being that if we are considering the F test then the upper hundred alpha percent points they are bigger than 1 there. So, therefore, a mod prudent thing would be to take the higher value that is S_2^2 by S_1^2 in case S_2^2 is bigger than S_1^2 , or S_1^2 by S_2^2 if S_1^2 is bigger than S_2^2 .

So, this follows F distribution on $m - 1, n - 1$ degrees of freedom, when $\sigma_1^2 = \sigma_2^2$. So, we can give a test based on this a convention is to place higher or you can say maximum of S_1^2 and S_2^2 in the numerator and the minimum in the denominator.

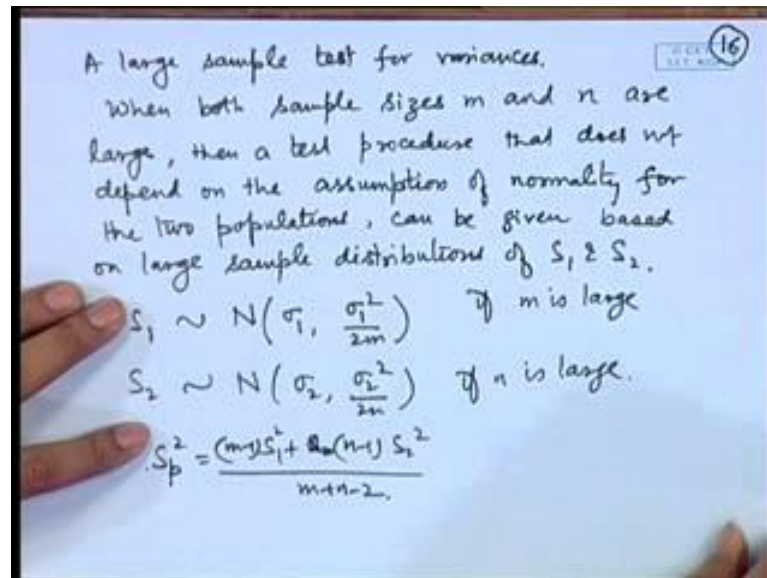
(Refer Slide Time: 16:32)



So, now if we are making a test of the hypothesis based on this F distribution on m minus 1, n minus 1 degrees of freedom then for 1 the test would be reject H_0 if F is greater than or equal to $f_{m-1, n-1, \alpha/2}$ or F less than or equal to $f_{m-1, n-1, 1-\alpha/2}$ of course, in the f distribution because of the reciprocal nature this is equal to 1 by $f_{m-1, n-1, \alpha/2}$. For the hypothesis 2, the rejection region will be reject H_0 if F is greater than or equal to $f_{m-1, n-1, \alpha}$. For the hypothesis testing problem 3 that is σ_1^2 less than σ_2^2 is alternative, we will be rejecting H_0 if F is less than or equal to $f_{m-1, n-1, 1-\alpha}$, which is equal to 1 by $f_{m-1, n-1, \alpha}$.

Once again in this situation also there may be a case when the normality assumption may not be satisfied; that means, the initial populations X_1 the X_1 sample and Y_1 sample etcetera they may not be from normal populations, in that case we may go for large sample approximations for the S_1^2 and S_2^2 .

(Refer Slide Time: 18:25)



And we may use the following a large sample test for variances, when both sample sizes m and n are large, then a test procedure that does not depend on the assumption of normality for the 2 populations, can be given based on large sample distributions of S_1 and S_2 . So, we can say that S_1 follows normal distribution with mean σ_1 , and variance, σ_1^2 by $2m$, if m is large and S_2 follows normal with mean σ_2 and variance σ_2^2 by $2n$, if n is large.

And once again we can utilize the fact that S_1 and S_2 are independently distributed, we may also create a pooled here that is S_p^2 that is equal to S_1 minus S_2 . Under the assumption that σ_1 is equal to σ_2 , we can consider like this that is m minus 1 S_1 square plus n minus 1, S_2 square divided by m plus n minus 2.

(Refer Slide Time: 20:53)

The image shows a handwritten slide with the following content:

$$Z^* = \frac{S_1 - S_2}{S_p \sqrt{\frac{1}{2m} + \frac{1}{2n}}} = \sqrt{\frac{2mn}{m+n}} \left(\frac{S_1 - S_2}{S_p} \right) \sim N(0,1)$$

So we can test for hypothesis about σ_1^2, σ_2^2 (①, ② & ③) as based on Z^* .

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 $Z^* \leq -z_\alpha$.

And we can construct the test statistic say Z star, that is equal to S 1 minus S 2 divided by Sp root 1 by 2 m, plus 1 by 2 m, which is actually equal to root 2 m n by m plus n, S 1 minus S 2 divided by Sp. So, this is approximately normal 0 1, because we are making the assumption that m and n are large.

So we can test for hypothesis about sigma 1 square and sigma 2 square, that is 1, 2 and 3 which I described earlier as based on Z star. So, for 1 we may say reject H naught, if modulus of Z star is greater than or equal to z alpha by 2. For the hypothesis testing problem 2, we will say reject H naught if Z star is greater than or equal to z alpha. For the third problem we will say reject H naught if Z star is less than or equal to minus z alpha. So, this test can be used when the assumption of normality for the basic populations is not very strong and of course, the sample sizes are large.

In the next lecture I will be discussing various problems where these tests can be utilized, we will also see a test for the proportions and for the binomial populations when we have 2 binomial populations.