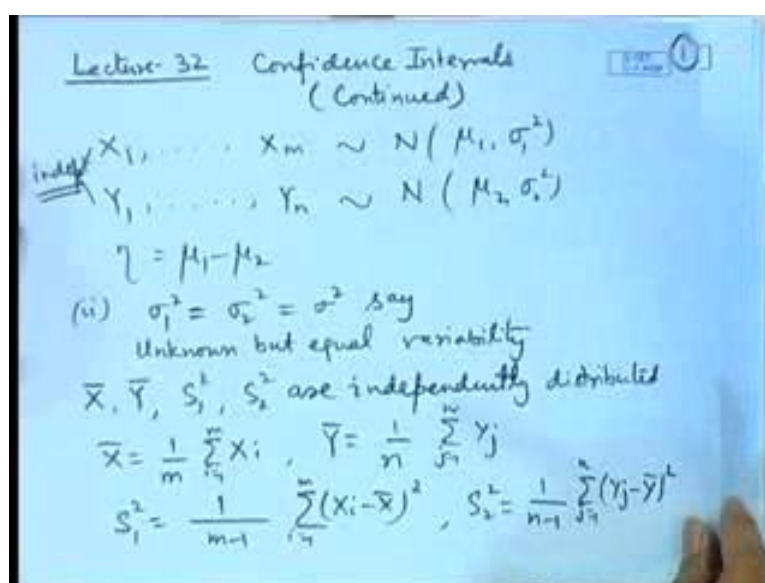


**Probability and Statistics**  
**Prof. Somesh Kumar**  
**Department of Mathematics**  
**Indian Institute of Technology, Kharagpur**

**Lecture – 63**  
**Confidence Intervals – III**

We continue our discussion on the Confidence Interval Estimation. Let me repeat the set up here.

(Refer Slide Time: 00:35)



We are interested in the comparison of the means of two normal populations. So, we have a sample  $X_1, X_2, \dots, X_m$  from normal  $\mu_1, \sigma_1^2$ ;  $Y_1, Y_2, \dots, Y_n$  is another independent random sample from normal  $\mu_2, \sigma_2^2$  population. These two samples are taken to be independent. So, we are interested in the confidence interval for  $\mu_1 - \mu_2$ , let us call it  $\eta$ . So, we have earlier found out the confidence interval for the situation when  $\sigma_1^2$  and  $\sigma_2^2$  are known. But in general the  $\sigma_1^2$  and  $\sigma_2^2$  may be unknown and we may be required to find out the confidence interval.

So, we take the case two that  $\sigma_1^2$  and  $\sigma_2^2$  are unknown but equal, that is unknown but equal variability. Now this type of situation may arise for example, you are looking at two brands of certain product. Now the variability of the say average life for example, it may be same, but average life themselves may be different. In such

cases this model is useful. Let us look at the analysis of this. So, as we have seen that the sampling distributions of  $\bar{X}$ ,  $\bar{Y}$ ,  $S_1^2$ ,  $S_2^2$  will be of interest here. So,  $\bar{X}$ ,  $\bar{Y}$ ,  $S_1^2$  and  $S_2^2$  are independent. In the sampling form normal distribution we know this fact independently distributed. Here  $\bar{X}$  is  $1$  by  $m$   $\sigma^2$   $\xi_i$  is equal to  $1$  to  $m$ ,  $\bar{Y}$  is the mean of the second sample that is  $1$  by  $n$   $\sigma^2$   $\eta_j$  is equal to  $1$  to  $n$ .

If we consider the sample variance of the first sample that is  $1$  by  $m-1$   $\sigma^2$   $\xi_i - \bar{X}$  whole square  $i$  is equal to  $1$  to  $m$  and  $S_2^2$  is equal to  $1$  by  $n-1$   $\sigma^2$   $\eta_j - \bar{Y}$  whole square that is the sample variance of the second sample.

(Refer Slide Time: 03:22)

Handwritten mathematical derivations on a blue background:

$$\bar{X} \sim N(\mu_1, \sigma^2/m)$$

$$\bar{Y} \sim N(\mu_2, \sigma^2/n)$$

$$\bar{X} - \bar{Y} \sim N\left(\frac{\mu_1 - \mu_2}{1}, \sigma^2\left(\frac{1}{m} + \frac{1}{n}\right)\right)$$

$$\frac{\bar{X} - \bar{Y} - \eta}{\sigma \sqrt{\frac{mn}{m+n}}} \sim N(0, 1)$$

$$\frac{(m-1)S_1^2}{\sigma^2} \sim \chi_{m-1}^2, \quad \frac{(n-1)S_2^2}{\sigma^2} \sim \chi_{n-1}^2$$

$$\frac{(m-1)S_1^2 + (n-1)S_2^2}{\sigma^2} \sim \chi_{m+n-2}^2$$

If we consider these quantities then we have the following observations; that is  $\bar{X}$  follows normal distribution with mean  $\mu_1$  and variance  $\sigma^2/m$ . So, here  $S_1^2$  and  $S_2^2$  both are same. Then  $\bar{Y}$  follows normal  $\mu_2$   $\sigma^2$  by  $n$ . If we consider here  $\bar{X} - \bar{Y}$  that will follow normal with mean  $\mu_1 - \mu_2$  and variance will be  $\sigma^2$   $1/m + 1/n$ .

So, if we want, this is the quantity  $\eta$ . So, we get  $\bar{X} - \bar{Y} - \eta$  divided by  $\sigma$  and root of this that is root of  $mn$  by  $m + n$  that will follow a standard normal distribution. However, this involves the unknown parameter  $\sigma$  also, so we cannot straight away use it as a pivot quantity. So, we need estimator for  $\sigma$  also. So, we can get it here by considering  $m-1$   $S_1^2$  by  $\sigma^2$  follows chi

square on  $m$  minus 1 degrees of freedom and  $n$  minus 1  $S_2$  is square by sigma square follows chi square on  $n$  minus 1 degrees of freedom.

Once again these two quantities are also independent. So, I can add these and we get  $m$  minus 1  $S_1$  square plus  $n$  minus 1  $S_2$  square divided by sigma square that follows chi square distribution on  $m$  plus  $n$  minus 2 degrees of freedom.

(Refer Slide Time: 05:17)

Let me define a quantity  $S_p$  square that is equal to  $m$  minus 1  $S_1$  square plus  $n$  minus 1  $S_2$  square divided by  $m$  plus  $n$  minus 2; that is pooled sample variance. If we use this pooled sample variance then what we are having is  $m$  plus  $n$  minus 2  $S_p$  square by sigma square is following chi square distribution on  $m$  plus  $n$  minus 2 degrees of freedom.

Now, we have the distribution of  $X$  bar minus  $Y$  bar minus  $\eta$  divided by sigma multiplied by a constant as a standard normal distribution, and let me call this quantity as say  $Z$  and I have a quantity let us call it say  $W$  this is having a chi square distribution. Another thing we can notice here this is that  $Z$  involving only  $X$  bar and  $Y$  bar and  $W$ 's involve in only  $S_1$  square and  $S_2$  square that is  $S_p$  square. So,  $Z$  and  $W$  are independently distributed.

So, if they are independently distributed I can look at the distribution of  $Z$  divided by  $W$  by  $m$  plus  $n$  minus 2 square root that will have  $t$  distribution on  $m$  plus  $n$  minus 2 degrees

of freedom. So, this quantity is equivalent to root mn by m plus n X bar minus Y bar minus eta divided by S p. So, that follows t distribution on m plus n minus 2 degrees of freedom.

Now let us observe: given the samples xi's and y j's we can evaluate X bar Y bar and S p and this involves the parameter eta for which we need the confidence interval and the distribution of this quantity is free from the parameters of the distribution. Therefore, this value t can be considered as a pivot quantity and we can make use of this to construct a confidence interval for eta that is mu 1 minus mu 2.

So, we look at the t distribution it is symmetric about the axis then it is symmetric about 0. So, this is f m plus n minus 2 t. So, we look at the point here this point is t alpha by 2 m plus n minus 2, and we have on the left hand side the single point that is minus t alpha by 2 m plus n minus 2. So, this intermediate probability is 1 minus alpha.

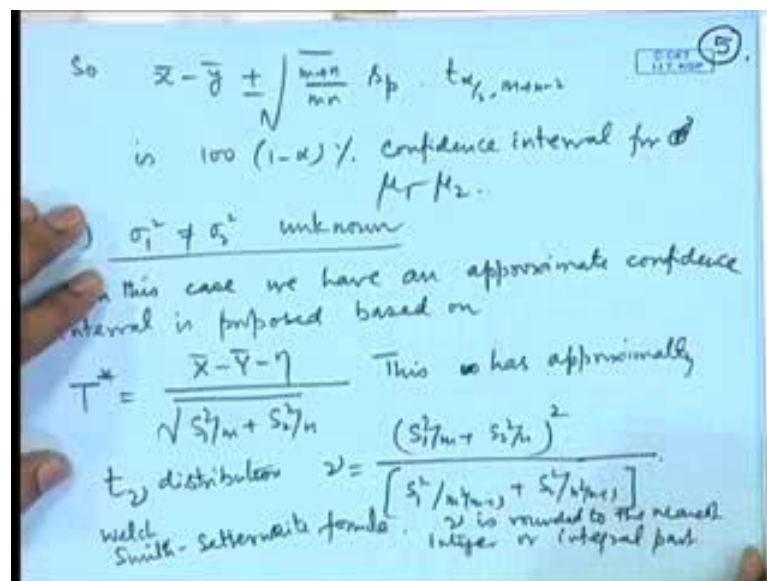
(Refer Slide Time: 08:37)

$$\begin{aligned}
 & P\left(-t_{\frac{\alpha}{2}, m+n-2} \leq T \leq t_{\frac{\alpha}{2}, m+n-2}\right) = 1 - \alpha \quad \text{④} \\
 \Leftrightarrow & P\left(-t_{\frac{\alpha}{2}, m+n-2} \leq \sqrt{\frac{mn}{m+n}} \frac{(\bar{X} - \bar{Y} - \eta)}{S_p} \leq t_{\frac{\alpha}{2}, m+n-2}\right) = 1 - \alpha \\
 \Leftrightarrow & P\left(-\sqrt{\frac{mn}{m+n}} S_p t_{\frac{\alpha}{2}, m+n-2} \leq \bar{X} - \bar{Y} - \eta \leq \sqrt{\frac{mn}{m+n}} S_p t_{\frac{\alpha}{2}, m+n-2}\right) = 1 - \alpha \\
 \Leftrightarrow & P\left(\bar{X} - \bar{Y} - \sqrt{\frac{mn}{m+n}} S_p t_{\frac{\alpha}{2}, m+n-2} \leq \eta \leq \bar{X} - \bar{Y} + \sqrt{\frac{mn}{m+n}} S_p t_{\frac{\alpha}{2}, m+n-2}\right)
 \end{aligned}$$

And we are in a position to write a statement that probability of minus t alpha by 2 m plus n minus 2 less than or equal to t is less than or equal to t alpha by 2 m plus n minus 2 that is equal to 1 minus alpha. So, expanding this t and then adjusting the terms we will be able to construct a confidence interval for mu 1 minus mu 2. So, t is here a square root of mn by m plus n X bar minus Y bar minus eta divided by S p; that is less than or equal to t alpha by 2 m plus n minus 2 that is equal to 1 minus alpha.

So, this is equivalent to  $\sqrt{m+n} \sqrt{\frac{mn}{m+n}} S_p t_{\alpha/2, m+n-2}$  less than or equal to  $\bar{X} - \bar{Y} - \eta$  less than or equal to  $\sqrt{m+n} \sqrt{\frac{mn}{m+n}} S_p t_{\alpha/2, m+n-2}$ ; that is equal to  $1 - \alpha$ . So, this means that  $\bar{X} - \bar{Y} - \sqrt{m+n} \sqrt{\frac{mn}{m+n}} S_p t_{\alpha/2, m+n-2}$  less than or equal to  $\bar{X} - \bar{Y} + \sqrt{m+n} \sqrt{\frac{mn}{m+n}} S_p t_{\alpha/2, m+n-2}$ ; that is equal to  $1 - \alpha$ .

(Refer Slide Time: 10:58)



So, in the situation when the variances of the two populations are unknown, but equal the confidence interval for  $\mu_1 - \mu_2$  is obtained as  $\bar{X} - \bar{Y} \pm \sqrt{m+n} \sqrt{\frac{mn}{m+n}} S_p t_{\alpha/2, m+n-2}$ . So, this is giving a  $100(1 - \alpha)$  percent confidence interval for  $\mu_1 - \mu_2$ .

Notice here is that since the variances were assumed to be equal we are making use of a pooled sample variance. Now one may ask a question that in place of this suppose we consider simply  $S_1^2$  or  $S_2^2$  only, because in that case also we are getting a variable which is having a distribution free from the parameters, so why not use only this or only this. The question is that if we use only say  $S_1^2$  then the degrees of freedom that we will get for the  $t$  variable will be  $m - 1$ . So, if we get only  $m - 1$  then in that case the interval will be having the width  $\bar{X} - \bar{Y} \pm \sqrt{m+n} \sqrt{\frac{mn}{m+n}} S_1 t_{\alpha/2, m-1}$ .

Now this term will not come here rather we will have  $S_1$  only this coefficient will not come here, here we will have only  $S_1$  and the degrees of freedom will be  $m - 1$ . Naturally the length of the interval will increase if we have less degree of freedom. So, in order to get more accuracy or you can say more precision we need a smaller interval with the same confidence coefficient. Therefore, it is beneficial to use more information here.

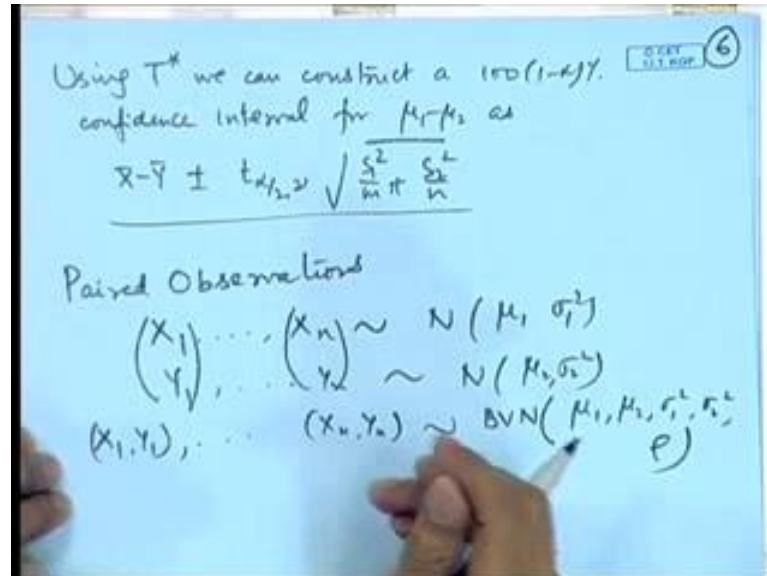
Let us take the case when both  $\mu_1$  and  $\mu_2$  may be unknown. Then let us look at the procedure here that has helped us to create this confidence interval. The procedure that we adopted was that the distribution of  $S_p^2$  by  $\sigma^2$  that is chi square and the Z variable that we utilized that also as a  $\sigma$  in the denominator. So, we were able to get rid of this. If the variances are not equal then in the first place we will be getting  $\sigma_1^2$  here and here we will get  $\sigma_2^2$ , so when we add the two terms in the denominators I will get  $S_1^2$  by  $\sigma_1^2$  and here  $S_2^2$  by  $\sigma_2^2$ . And the same thing will happen with the Z also, where we will get  $\sigma_1$  square by  $m$  plus  $\sigma_2$  square by  $n$ .

So, in no way by taking the ratios I can get rid of  $\sigma_1^2$  and  $\sigma_2^2$ . Actually it turns out that there is no exact confidence interval; that means, the interval which is having the length a shortest length and as well as a fix confidence coefficient; that means, a distribution free term here not getting. In this case this is known as a Behrens-Fisher situation. So, we will consider this case  $\sigma_1^2$  is not equal to  $\sigma_2^2$  and unknown; that means, they are completely unknown. In this case a approximate an approximate confidence interval is proposed based on; let us call it  $t^*$  that is  $\bar{X} - \bar{Y} - \eta$  divided by square root  $S_1^2$  by  $m$  plus  $S_2^2$  square by  $n$ .

So, how this has come? In the first case where  $\sigma_1^2$  by  $m$  plus  $\sigma_2^2$  is square by  $n$  was there we have simply replaced  $\sigma_1^2$  and  $\sigma_2^2$  square by their unbiased estimates. So, it was proved by Welch etcetera that this is having has approximately t distribution on  $\nu$  degrees of freedom, where  $\nu$  is given by  $S_1^2$  square by  $m$  plus  $S_2^2$  square by  $n$  whole square divided by  $S_1^2$  square by  $m$  square into  $m - 1$  plus  $S_2^2$  square by  $n$  square into  $n - 1$ ; which is by Welch and it is known as a Smith Satterthwaite formula.

Now this need not be an integer. So, nu is rounded off to the nearest integer or integral part. That means, suppose it is turning out to be 11.3 second we take only 11.

(Refer Slide Time: 17:04)



So, using this one can write a confidence interval; using T star we can construct a 100 1 minus alpha percent confidence interval for mu 1 minus mu 2 as X bar minus Y bar plus minus t alpha by 2 mu a square root of S 1 square by m plus S 2 square by n. This will be the confidence interval. Then there is no information about the equality of sigma 1 square and sigma 2 square.

Now there is another situation which occurs quite frequently; for example, we are considering the comparison of the two training procedures. So, suppose there are two training procedures for certain learning. So, we select say 10 pupils and we give them instructions using one training procedure a test is conducted to measure the outcome of that. Now, for the same set of 10 pupils another learning procedure is important for a fix period of time and another test is conducted.

Now the scores are not independent because our subjects are not independent, same set of people has been selected. For example, it could be some weight reduction procedures like the fatty people are there and we are giving them certain weight reduction program. So, by taking certain procedure for one month their weight is reducing by this much. Now for the same set of people another procedure is adopted then how much weights have been reduced. So, we compare the same set of people with respect to their course.

So, here this is related to paired observations. So, here although you are saying  $X_1, X_2, \dots, X_n$  say follow normal  $\mu_1, \sigma_1^2$  and  $Y_1, Y_2, \dots, Y_n$  they follow normal  $\mu_2, \sigma_2^2$ , but actually the sample has not been selected in this way because these observations may be paired. So, basically the model becomes that  $X_1, Y_1, X_2, Y_2, \dots, X_n, Y_n$  this is having some sort of bivariate normal distribution with parameters  $\mu_1, \mu_2, \sigma_1^2, \sigma_2^2$  and some correlation coefficient  $\rho$  may be there. Once again we are interested in the interval for  $\mu_1 - \mu_2$  that is we want to look at the difference in the average effectiveness etcetera.

A simple procedure for this is obtained by using the linearity property of bivariate normal distribution, because we know that if the random variable  $X, Y$  is having a bivariate normal distribution then any linear combination  $aX + bY$  is again having a univariate normal distribution.

(Refer Slide Time: 20:33)

confidence interval for  $\mu_1 - \mu_2$  is

$$\bar{X} - \bar{Y} \pm t_{\alpha/2, n-1} \sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{n}}$$

Paired Observations

$$(X_1), \dots, (X_n) \sim N(\mu_1, \sigma_1^2)$$

$$(Y_1), \dots, (Y_n) \sim N(\mu_2, \sigma_2^2)$$

$$(X_1, Y_1), \dots, (X_n, Y_n) \sim \text{BVN}(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$$

$$d_i = X_i - Y_i \sim N(\underbrace{\mu_1 - \mu_2}_{\mu_d}, \underbrace{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}_{\sigma_d^2})$$

So, here if I make use of say observations let me call it  $d_i$  that is equal to  $x_i - y_i$ , then that will follow normal distribution with mean  $\mu_1 - \mu_2$  and some variance let me call it  $\sigma_d^2$ ; actually it will be  $\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2$ . So,  $\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2$ , will let me call it  $\sigma_d^2$ . That is not important here because they are all unknown and we need only an estimate of this, because we are interested here in the confidence interval about  $\mu_1 - \mu_2$ .



So we can make use of; now this looks like a problem of the confidence interval for a mean of a normal distribution which we have done in the first place.

(Refer Slide Time: 21:30)

Handwritten mathematical derivations on a blue background:

$$\bar{d} = \frac{1}{n} \sum_{i=1}^n d_i, \quad S_d^2 = \frac{1}{n-1} \sum (d_i - \bar{d})^2$$

$$\bar{d} \sim N(\eta, \sigma_d^2/n), \quad \frac{\sqrt{n}(\bar{d} - \eta)}{\sigma_d} \sim N(0, 1)$$

independent

$$\frac{(n-1)S_d^2}{\sigma_d^2} \sim \chi_{n-1}^2$$

$$\frac{\sqrt{n}(\bar{d} - \eta)}{S_d} \sim t_{n-1}$$

A normal distribution curve is shown with a central peak and two vertical lines at  $-t_{\alpha/2, n-1}$  and  $t_{\alpha/2, n-1}$ . The area between these lines is labeled  $1 - \alpha$ .

$\left( \bar{d} - \frac{S_d}{\sqrt{n}} t_{\alpha/2, n-1}, \bar{d} + \frac{S_d}{\sqrt{n}} t_{\alpha/2, n-1} \right)$  is a  $100(1 - \alpha)\%$  confidence interval for  $\eta = \mu_1 - \mu_2$ .

So, we can consider say  $\bar{d}$  as  $\frac{1}{n} \sum_{i=1}^n d_i$  is equal to  $\frac{1}{n} \sum_{i=1}^n d_i$ , and we consider  $S_d^2$  as  $\frac{1}{n-1} \sum_{i=1}^n (d_i - \bar{d})^2$ . Say if we look at this then we can see the  $\bar{d}$  follows normal  $\eta, \sigma_d^2/n$  and from here we can get  $\bar{d} - \eta$  by  $\sigma_d/\sqrt{n}$  follows normal  $0, 1$ . Also  $(n-1)S_d^2/\sigma_d^2$  will follow chi square distribution on  $n-1$  degree of freedom. And once again these two variables will be a statistically independent. So, using this we can write a square root  $n(\bar{d} - \eta)/S_d$  that will be having a  $t$  distribution on  $n-1$  degree of freedom.

Now, observe this function here it is involving the random variables that are observations  $x_i$  and  $y_i$ 's,  $\bar{d}$  is the mean calculated from the differences and  $S_d^2$  is calculated as the variance of the difference observations. And here the parameter of interest  $\eta$  is appearing and  $\sigma_e^2$  etcetera are absent here. So, this can be used as a pivot quantity and we get a confidence interval by writing down from the distribution of the  $t$  on  $n-1$  degree of freedom. So, this probability is  $1 - \alpha$  and we get  $\bar{d} - \frac{S_d}{\sqrt{n}} t_{\alpha/2, n-1}$  to  $\bar{d} + \frac{S_d}{\sqrt{n}} t_{\alpha/2, n-1}$ . So, this becomes  $100(1 - \alpha)\%$  confidence interval for  $\eta$  that is equal to  $\mu_1 - \mu_2$ .

So, we observe here that all these cases are differently handled; that is when we observe a sample we have to look at carefully. So, if the variance is known to us then we have some procedure, if the variances are unknown but we suspect that the variances may be equal then we have another procedure, if the variances are completely unknown then we have another procedure. On the other hand if the sampling is not done in the independent fashion; that means, we have correlated observations then we have may arrange the data in a paired way and then we can apply a pairing formula.

So, the confidence interval for the same parameter  $\mu_1$  minus  $\mu_2$  when we are sampling from to normal populations it is dependent upon the situation. We have to a statistician has to carefully see that which type of method will be adopted here for finding out the confidence interval, otherwise you will be coming up with the faulty conclusions.