**Module No. #01**
**Lecture No. #32**
**Estimation-VI**

(Refer Slide Time: 00:24)



So, we continue our discussion on the confidence interval estimation. Let me repeat the setup here, we are interested in the comparison of the means of two normal populations. So, we have a sample X1, X2, Xm from normal mu1, sigma1 square, Y1, Y2, YN is another independent random sample from normal mu2, sigma2 square population- these two samples are taken to be independent. So, we are interested in the confidence interval for mu1 minus mu2, let us call it eta. So, we have earlier found out the confidence interval for the situation when sigma1 square and sigma2 square are known, but in general, the sigma1 square and sigm 2 square may be unknown and we may be required to find out the confidence interval.

So, we take the case two, that sigma1 square and sigma2 square are unknown, but are 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. So, now the variability of the

say, average life for example, it may be same, but average lives themselves may be different, so 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 X bar, Y bar S1 square S2 square will be of interest here. So, X bar, Y bar, S1 square and S2 square are independent- in the sampling from normal distribution we know this fact, independently distributed.

Here X bar is 1 by m sigma Xi, i is equal to 1 to m; Y bar is the mean of the second sample that is, 1 by n sigma yj, j is equal to 1 to n. If we consider the sample variance of the first sample that is, 1 by m minus 1 sigma Xi minus X bar whole square, i is equal to 1 to m and S2 square is equal to 1 by n minus 1 sigma Yj minus Y bar whole square, that is the sample variance of the second sample.

(Refer Slide Time: 03:19)



If we consider these quantities, then we have the following observation: that is, X bar follows normal distribution with mean mu1 and variance sigma square by m, so, here sigma1 square and sigma2 squares both are same. Then, Y bar follows normal mu2 sigma square by n. So, if we consider here, X bar minus Y bar, that will follow normal with mean mu1 minus mu2 and variance will be sigma square 1 by m plus 1 by n.

So, if we want, so, this is the quantity eta. So, we get X bar minus Y bar minus eta divided by sigma and root of this that is, root of mn by m plus 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 an estimator for sigma also. So, we can get it here by considering m minus 1 S1 square by sigma square follows chi square on m minus 1 degrees of freedom, and n minus 1 S2 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 this and we get m minus 1 S1 square plus minus 1 S2 square divided by sigma square, that follows chi square distribution on m plus n minus 2 degrees of freedom.

(Refer Slide Time: 05:15)



Let me define a quantity Sp square, that is equal to m minus 1 S1 square plus n minus 1 S2 square divided by m plus n minus 2- that is, pooled sample variance. If we used this pooled sample variance, then what we are having is m plus n minus 2 Sp 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 is that Z is involving only X bar and Y bar, and W is involving only S1 square and S2 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 Sp, so, that follows t distribution on m plus n minus 2 degrees of freedom.

Now, let us observe, given the samples Xis and Yjs, we can evaluate X bar, Y bar and Sp, 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, mu1 minus mu2.

So, we look at the t distribution, it is symmetric about the axis, it is symmetric about zero and, so, this is fm 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 similar point, that is minus t alpha by 2 m plus n minus 2. So, this intermediate probability is 1 minus alpha.

(Refer Slide Time: 08:34)



$$P\left(-t_{\frac{\alpha}{2},\,m+n-2} \leq T \leq t_{\frac{\alpha}{2},\,m+n-2}\right) = 1-\alpha$$

$$\Leftrightarrow P\left(-t_{\frac{\alpha}{2},\,m+n-2} \leq \sqrt{\frac{mn}{m+n}}\,\frac{(\bar{X}-\bar{Y}-\eta)}{Sp} \leq t_{\frac{\alpha}{2},\,m+n-2}\right) = 1-\alpha$$

$$P\left(-\sqrt{\frac{m+n}{mn}}\,Sp\,t_{\frac{\alpha}{2},\,m+n-2} \leq \bar{X}-\bar{Y}-\eta \leq \sqrt{\frac{m+n}{mn}}\,Sp\,t_{\frac{\alpha}{2},\,m+n-2}\right) = 1-\alpha$$

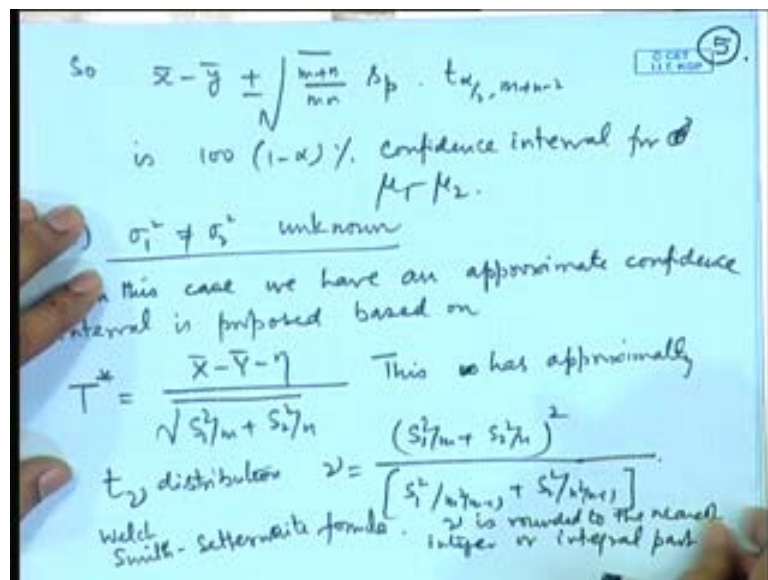$$\Leftrightarrow P\left(\bar{X}-\bar{Y}-\sqrt{\frac{m+n}{mn}}\,Sp\,t_{\frac{\alpha}{2},\,m+n-2} \leq \eta \leq \bar{X}-\bar{Y}+\sqrt{\frac{m+n}{mn}}\,Sp\,t_{\frac{\alpha}{2},\,m+n-2}\right) = 1-\alpha$$

And we are in a position to write the 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 mu1 minus mu2. So, T is here square

root of m n by m plus n X bar minus Y bar minus eta divided by Sp, 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 root m plus n by mn Sp t alpha by 2 m plus n minus 2 less than or equal to X bar minus Y bar minus eta less than or equal to square root m plus n by mn Sp t alpha by to m plus n minus 2, that is equal to 1 minus alpha.

So, this means that X bar minus Y bar minus root m plus n by mn Sp t alpha by 2 m plus n minus 2 less than or equal to X bar, less than or equal to eta less than or equal to X bar minus Y bar plus, that is equal to 1 minus alpha.

(Refer Slide Time: 10:54)



So, in the situation when the variances of the two populations are unknown, but equal, the confidence interval for mu1 minus mu2 is obtained as X bar minus Y bar plus minus square root m plus n by mn Sp t alpha by 2 m plus n minus 2, so, this is giving a 100 1 minus alpha percent confidence intervals for mu1 minus mu2.

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 S1 square or S2 square only, because in that case also we are getting a variable which is free from the, 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, S1 square,

then the degrees of freedom that we will get for the T variable will be m minus 1, so, if we get only m minus 1, then in that case, the interval will be having the width X bar minus Y bar plus minus square root m plus n by mn, now, this term will not come here, rather we will have S1 only, this coefficient will not come here, here we will have only S1, and the degrees of freedom will be m minus 1 naturally, the length of the interval will increase if we have less number of, less degrees 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 mu1 and mu2 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 Sp square by sigma square that is, chi square and the Z variable that we utilized, that also has 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 sigma1 square here and here we will get sigma2 square, so when we add the two terms in the denominators I will get S1 square by sigma1 square and here S2 square by sigma2 square, and the same thing will happen with the Z also, where we will get sigma1 square by m plus sigma2 squares by n. So, in no way by taking the ratios I can get rid of sigma1 square and sigma2 square actually, 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 fixed confidence coefficient that means, the distribution free term we are not getting.

In this case, this is known as a variance special situation. So, we will consider this case, sigma1 square is not equal to sigma2 squares 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 star, that is X bar minus Y bar minus eta divided by square root S1 square by m plus S2 square by n. So, how this has come? In the first case where sigma1 square by m plus sigma2 squares by n was there, we have simply replaced sigma1 square and sigma2 squares by their unbiased estimates. So, it was proved by Welch, etcetera, that this is having, has approximate, this has approximately t distribution on nu degrees of freedom, where nu is given by S1 square by m plus S2 square by n whole square divided by S1 square by m square into m minus 1 plus S2 square by n square into n minus 1- it is by Welch and it is known as smith-satterthwaite formula.

So, 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 11.37 we take only 11.

(Refer Slide Time: 17:02)



So, using this one can write a confidence interval. Using T star we can construct a 100 1 minus alpha percent confidence interval for mu1 minus mu2 as X bar minus Y bar plus minus t alpha by 2 nu square root of S1 square by m plus S2 square by n- this will be the confidence interval when there is no information about the equality of sigma1 square and sigma2 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, ten peoples 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 ten peoples another learning procedure is imparted for a fixed period of time and another test is conducted. Now, the scores are not independent because our subjects are not independent, same people, same set of people has been selected. For example, it could be some weight reduction procedure 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 scores.

So, here this is related to paired observations, paired observations. So, here although you are saying X1, X2, XN say, follow normal mu1 sigma1 square and Y1, Y2, YN, they follow normal mu2 sigma2 squares, 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 Y1, X2 Y 2, Xn Yn, this is having some sort of bivariate normal distribution with parameters mu1 mu2 sigma1 square sigma2 square and some correlation coefficient, rho may be there.

Once again we are interested in the interval for mu1 minus mu2 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 plus bY is again having a univariate normal distribution, so, here if I make use of say, observations, let me call it di, that is equal to Xi minus Yi, then that will follow normal distribution with mean mu1 minus mu2 and some variance, let me call it sigma D square- actually it, will be sigma1 square plus sigma2 square minus twice rho sigma1 sigma2- so, sigma1 square plus sigma2 square minus twice rho sigma1 sigma2, let me call it sigma D square, 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 mu1 minus mu2.

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: 21:27)



So, we can consider say d bar as 1 by n sigma di, i is equal to 1 to n, and we consider Sd square as 1 by n minus 1 sigma di minus d bar whole square. So, if you look at this, then we can see that d bar follows normal eta sigma d square by n, and, so, from here we can get d bar minus eta root n by sigma d follows normal 0, 1; also n minus 1 Sd square by sigma square by sigma d square, that will follow chi square distribution on n minus 1 degrees of freedom, and once again, these two variables will be statistically independent. So, using this we can write square root n d bar minus eta divided by Sd, that will be having a t distribution on n minus 1 degrees of freedom.

Now, observe this function here, it is involving the random variables, that is observations Xis and Yis, d bars are the mean calculated from the differences and Sd square is calculated as the variance of the difference observations, and here the parameter of interest eta is appearing and sigma e squares, 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 minus 1 degrees of freedom. So, this probability is 1 minus alpha and we get d bar minus sd by root n t alpha by 2 n minus 1 to d bar plus sd by root n t alpha by 2 n minus 1, so this becomes 100 1 minus alpha percent confidence intervals for eta, that is equal to mu1 minus 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 may arrange the data in a paired way and then we can apply a pairing formula. So, the confidence interval for the same parameter mu1 minus mu2 when we are sampling from two 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, he will be coming up with the faulty conclusions.

(Refer Slide Time: 25:15)



Let me take up some examples here to illustrate the situations. So, to compare, to compare the strength, the gripping strengths of left hand and right hand of ten left handed, of left handed persons, the measurements are made on ten persons and the data is observed. So, left hand and right hand, and we have persons 1, 2, 3, 4, 5, 6, 7, 8 9 and 10; the gripping strengths are measured as 140, 90, 125, 130, 95, 121, 85, 97, 131, 110; for the right hand it is 138, 87, 110, 132, 96, 120, 86, 90, 129, 100.

So, we need the confidence interval for say mu1 minus mu2. Now, observe here that this is the data related with the correlated observations, so, we will need here the means of, so, let me call this as the first set, so, this is Xi data this is Yi data. So, we will look at dis, the differences here; so, the differences here is 2, 3, 15, minus 2, minus 1, 1, minus 1, 2, and 10. So, we look at the d bar value here, which is the mean of this that is, 20... so, 17... 24... 26... 36... so, that is 3.6. Similarly, we calculate sd, that will be equal to 1 by 9 sigma di square minus d bar square. So, once again, it can be easily evaluated it is 4 plus 9 plus 225 plus 4 plus 1 plus 1 plus 1 plus 49 plus 4 plus 100 minus 3.6 square- so, this value can be evaluated.

Now, we look at the value of t on, suppose we want a 90 percent confidence interval, so we need .059 that is equal to 1.83. So, we get the confidence interval as 3 .6 plus minus sd by root 10 into 1.833, that will be the confidence interval for the difference in the gripping strengths of left hand and the right hand of the left handed persons.

(Refer Slide Time: 29:44)



Let me take another example here, to compare age at marriage of women in two ethnic groups a random sample of hundred women is taken, and we observe that X bar is equal to 18.5 years, Y bar is equal to 20.7 years and S1 is equal to 5.8, S2 is equal to 6.3 and we want say, a confidence interval for this. So, we calculate here that we may use the model for Sp square. So, S p square is equal to m minus 1 S1 square plus n minus 1 S2

square divided by m plus n minus 2, that is equal to 99 S1 square plus 99 S2 square by 198, that is equal to 5.8 square plus 6. 3 square by 2- so, this value can be evaluated.

In a similar way, so, we have the confidence interval as X bar minus Y bar minus root m plus n by mn that is, and then, Sp into t alpha by 2- suppose I want 90 percent confidence interval- so, 0.05 and the degrees of freedom will be m plus n minus 2; so, this value we can see t 0.05, 198 which is almost as a normal distribution 1645. So, we substitute these values here, 18.5 minus 20. 7 minus- this is 100 plus 100, that is 200 by 100- so, that is root 2 by 10 into 1.645, so, plus minus. That gives the confidence interval for the difference in the ages at marriage of women in two ethnic groups. So, here we have the pooled formula, we may actually do a testing of hypothesis for sigma1 square is equal to sigma2 square, and if sigma1 square is equal to sigma2 square is accepted, then we may go for this formula.
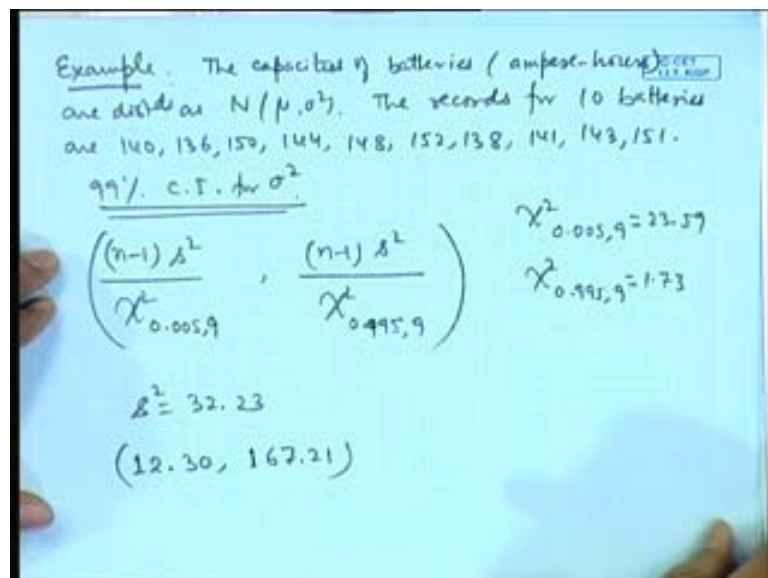
(Refer Slide Time: 33:00)



Let me take another example here. Two machines are used to fill plastic bottles with dishwashing detergent. The standard deviations of fill volume are known to be sigma1 is equal to 0.1 5 fluid ounces and sigma2 is equal to 0.12 fluid ounces for the two machines. Now, two random samples of n1 is equal to 12 bottles from machine one and n2 is equal to ten bottles from machine two are selected and the observations are X1 bar

is equal to 30.87, X 2 bar is equal to 30.68. So, find 90 percent confidence interval for mu1 minus mu2.

So, here we can see, we can look at the confidence interval as X bar minus Y bar plus minus square root sigma1 square by m plus sigma2 square by n z0.05- now, z0.05, we can see from the tables of normal distribution, it is 1. 645. So, this interval becomes 30.87 minus 30.6 8 plus minus square root, now, sigma1 square is 0.15 square by 12 plus, sigma2 square is 0.12 square by n, n is 10 multiplied by 1.645 So, after simplification, these values turns out to be 0.095 to 0.285, so, this is 90 percent confidence interval for the main difference that, is mu1 minus mu2. So, here the variances were known, sigma1 square and sigma2 square, so we have adopted a procedure where the formula for non variances is utilized here.

(Refer Slide Time: 36:53)



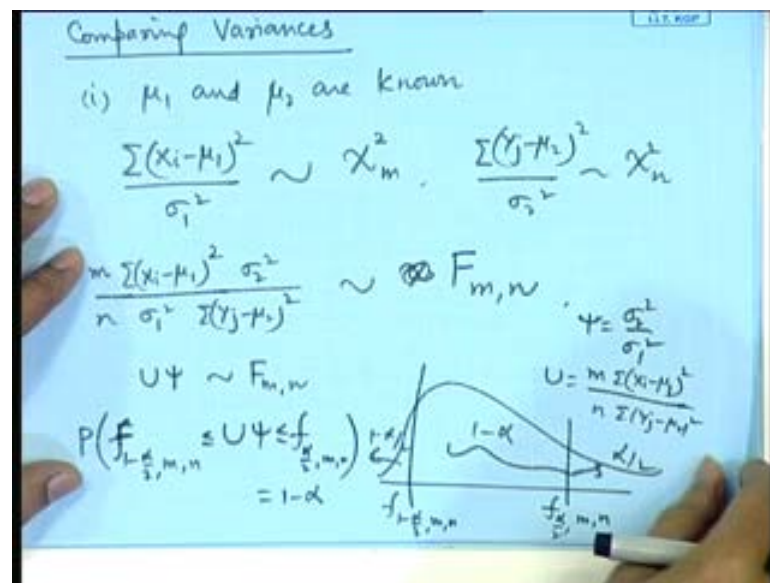Let me take another example here. The capacities of batteries, so, these are measured in say ampere-hours; they are distributed as normal mu, sigma square. The records for ten batteries are say 140, 136, 150, 144, 148, 152, 138, 141, 143, 151. We want 99 percent confidence interval for sigma square.

So, now, here, we will make use of the fact that mu is unknown. So, if mu is unknown, then the formula for confidence interval for sigma square is based on chi square on n

minus 1 degrees of freedom, the formula is n minus 1 s square by chi square. So, 0.005 n is 1, so this is 9 to n minus 1 s square by chi square 0.0959, s<mark>o, 995.</mark> So, these values we see from the tables of the chi square distribution that is, chi square <mark>0.0059</mark>, it is 23.99 and chi square 0.995 on 9 degrees of freedom is 1.73. So, s square we calculate here, it is turning out to be 32.23. So, after substitution of these values, the confidence interval turns out to be 12.30 to 1 67.21, which is pretty large confidence interval, but that will be there because we are considering for sigma square and the variability of the original sample itself is large, this is, s square is 32.23 here. If we reduce the confidence level, suppose we make it 90 percent, then this will be shrinking, since we have made a very high confidence level that is why the confidence interval is very large, which looks slightly impractical also.

(Refer Slide Time: 40:13)



Next, we look at the confidence intervals for variances- so, comparing variances. Again, we have two cases, that is mu1 and mu2 are known; if mu1 and mu2 are known, then we make use of sigma Xi minus mu1 square by sigma1 square following chi square distribution on m degrees of freedom, and sigma Yj minus mu2 square by sigma2 square follows chi square distribution on n degrees of freedom. So, if you take the ratios here, sigma Xi minus mu1 square by sigma1 square (( )) m, so, that is m here, divided by sigma Yj minus mu2 squares by sigma2 square, so, that will come in the numerator, divided by n, that will have chi square, f distribution on m and n degrees of freedom.

So, if we look at this quantity, if mu1 and mu2 are known, then here the ratio sigma2 square by sigma1 square is coming, let us denote it by say, psi that is sigma2 square by sigma1 square; so, we are having, and let me use the notation say, U as m sigma Xi minus mu1 square divided by n sigma Yj minus mu2 square. So, if you look at this one, then we are having U psi following f distribution on m n degrees of freedom. So, if we make use of the tables of f distribution that is, f on m and n degrees of freedom here, and f1 minus alpha by 2 on m n degrees of freedom, this is alpha by 2 and this is 1 minus alpha by 2, so, this is 1 minus alpha; so, probability of f 1 minus alpha by 2 m n less than or equal to U psi less than or equal to f alpha by 2 m n, that is equal to 1 minus alpha.

(Refer Slide Time: 42:58)



So, we can write probability of U, so, divided by f of alpha by 2 m n less than or equal to sigma1 square that is, 1 by psi, it becomes sigma1 square by sigma2 square less than or equal to U divided by f 1 minus alpha by 2 m n, that is equal to 1 minus alpha. So, we have a 1 minus alpha confidence interval for sigma1 square by sigma2 square, this can also be written as U f 1 minus alpha by 2 n m to U f alpha by 2 n m, by using the ratio, or you can say reciprocal property of the f distribution, because we know that 1 by f m n is 1 by, is equal to f n m, so, this property can be utilized here.

(Refer Slide Time: 45:14)



Let me give one example here for confidence interval for the ratios. So, two brands of say, cough medicine are given and the response times are measured in days. So, here we are having the data, m is equal to say 10, n is equal to 12, S1 and, we are getting the observations as x1 is equal to say, 2, 3, 2, 4, 2, 5, 6, so, 3, 7 and then, 1, 2, so, we have ten data here and for y we have the data say, 3, 4, 6, 8, 3, 2, 9, 5, 11, 7, 2, 1. Now, based on this we calculate x bar y bar and we calculate sigma xi minus mu1. So, it is given that mu1 is say, 3 and mu2 is equal to 5. So, if we are looking at sigma xi minus mu1 square and sigma yj minus mu2 square, then that will follow chi square on 9 and this divided by sigma2 square follow chi square on 11 degrees of freedom. So, we can construct ten sigma xi minus mu1 square by 12 sigma yj minus mu2 square, and then, we need to look at the tables of f on say 0.5, 10 and 12 degrees of freedom.

Now, another situation may occur when mu1 and mu2 are unknown. If mu1 and mu2 are unknown, then we will not be able to make use of the formula that we derived earlier because there in the confidence interval mu1 and mu2 are actually appearing, so what we do, we make use of S1 square and S2 square. So, we have m minus 1 S1 square follows chi square distribution on m minus 1 degrees of freedom and n minus 1 S2 square by sigma2 square follows chi square distribution on n minus 1 degrees of freedom. Furthermore, these two random variables are independent. So, we can make use of the ratios m minus 1 S1 square by sigma1 square divided by m minus 1 divided by n minus 1 S2 square by sigma2 square into n minus 1, that will follow f distribution on m minus 1, n minus 1 degrees of freedom, which is reducing to sigma2 square by sigma1 square S1 square by S2 square, this follows f distribution on m minus 1 n minus 1 degrees of freedom.

(Refer Slide Time: 49:15)



So, making use of distribution of f that is, we have f alpha by 2 m minus 1 n minus 1 and f 1 minus alpha by 2 m minus 1 n minus 1, intermediate probability is 1 minus alpha; so, probability that f 1 minus alpha by 2 m minus 1 n minus 1 is less than or equal to sigma2 square by sigma1 square S1 square by S2 square is less than or equal to f alpha by 2 m minus 1 n minus 1, that is equal to 1 minus alpha. So, we make use of this and adjust the coefficients as probability that S2 square by S1 square f 1 minus alpha by 2 m minus 1 n minus 1 less than or equal to sigma2 square by sigma1 square less than or equal to S2 square by S1 square f alpha by 2 m minus 1 n minus, 1 that is equal to 1 minus alpha. So, we are getting 100 1 minus alpha percent confidence intervals for sigma2 square by sigma1 square. We can reverse it, if we want for sigma1 square by sigma2 square, then we interchange the roles here, we put S1 square by S2 square and the degrees of freedom will get reverse, it will become n minus 1 m minus 1.

(Refer Slide Time: 50:48)



So, we give one example here. So, S2 square by S1 square f 1 minus alpha by 2 m minus 1 n minus 1 to S2 square by S 1 square f alpha by 2 m minus 1 n minus 1 is 100 1 minus alpha percent confidence interval for sigma2 square by sigma1 square.

(Refer Slide Time: 51:11)



So, say viscosity of two brands of oil used in cars is measured and the following data is recorded. So, from brand one you have 10.62, 10.58, 10.33, 10.72, 10.44. For brand two

it is 10.50, 10.52, 10.62, 10.53. Suppose we want a confidence interval for sigma2 squares by sigma1 square. So, we will calculate the values here, s1 square s2 square; so, s1 square turns out to be 0.02362, s2 square is equal to 0.002825, you can see here there is a 10 times difference here. So, the f values that 22 square by, or you can say s1 square by s2 square will be equal to 8.36. So, if you look at the f value on 0. 5 say, 1, 2, 3, 4, 5, so, 4, 3 degrees of freedom, that is equal to 9.1172, and f value 0.9543, that is equal to 0.1517.

So, a 90 percent confidence interval for sigma1 square by sigma2 square, that will be equal to 8.36 into 0.1517 to 8.36 into 9.1172. So, this is the confidence interval for the ratio of the variances here.

So, in a given practical situation we need to analyze that what is the model that will be applicable and accordingly we make use of the formulae. So, for example, when we are looking at the confidence intervals for mu1 minus mu2, then we worry about that what is the status of the variances, if the variances are known, then we have some formula, if the, that is, based on the Z that is, normal distribution, if we have variances unknown, but equal, then we have a formula which is based on T distribution, based on the pooling of the concept, pooling of the variances, and if we have variances to be completely unknown, then in that case we have another approximate T distribution formula and we make use of that.

On the other hand, if the data is correlated, then we make use of pairing and a paired T formula is used. Similarly, when we are worried about the confidence interval for the sigma1 square and sigma2 square, then we look at the knowledge about the means, if the means are known, then we have a formula based on f distribution on the total degrees of freedom m and n, if the means are unknown, then we have another formula which is based on S1 square S2 square and the degrees of freedom are slightly reduced to m minus 1, n minus 1.

Now, these formulae are quite standard because they are making use of the sampling distribution from the normal populations. When we do not have normal populations then in that case, we may have to look for appropriate sampling distribution, for example, if we are dealing with uniform distribution, if we are dealing with from the exponential

populations, then we look at from the description that what is the sufficient statistics, from there we find out the pivoting quantity if we are able to derive the sampling distribution of that. So, the techniques for that and also for the proportions are available and one can work out various formulae for confidence interval from other populations as well. So, that is part of another course, that is statistical inference that, we will be doing later on.