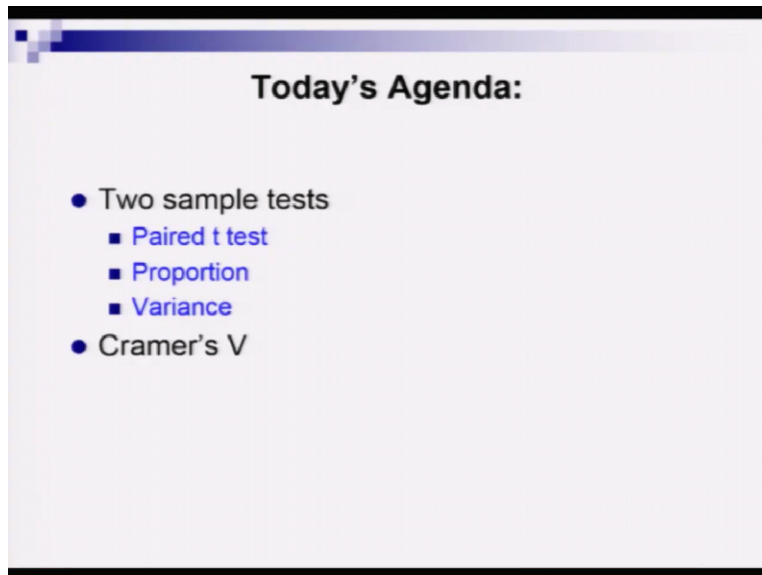


Applied Statistics and Econometrics
Professor Deep Mukherjee
Department of Economic Sciences
Indian Institute of Technology, Kanpur
Lecture-14
Hypothesis Testing (Part IV)

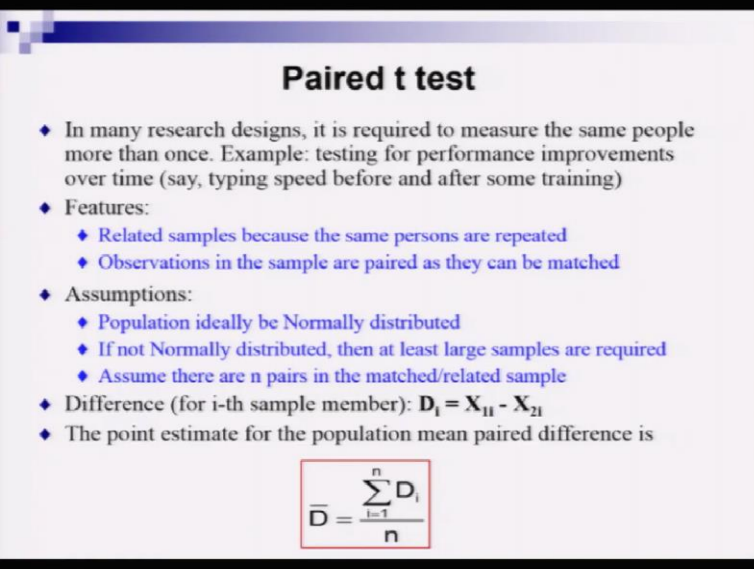
Hello, friends. Welcome back to the lecture series on Applied Statistics and Econometrics. So today, we are going to continue our discussion on hypothesis testing and we are going to look at a couple of new tests. So, let us have a look at today's agenda item.

(Refer Slide Time: 00:31)



So, in today's lecture, we are going to talk about a very useful test called paired t test. Then we are also going to look at the other variations of 2 sample tests, namely proportion and variance tests. And finally, we will end our discussion by briefly discussing a very interesting concept called Cramer's V.

(Refer Slide Time: 01:38)



Paired t test

- ◆ In many research designs, it is required to measure the same people more than once. Example: testing for performance improvements over time (say, typing speed before and after some training)
- ◆ Features:
 - ◆ Related samples because the same persons are repeated
 - ◆ Observations in the sample are paired as they can be matched
- ◆ Assumptions:
 - ◆ Population ideally be Normally distributed
 - ◆ If not Normally distributed, then at least large samples are required
 - ◆ Assume there are n pairs in the matched/related sample
- ◆ Difference (for i -th sample member): $D_i = X_{1i} - X_{2i}$
- ◆ The point estimate for the population mean paired difference is

$$\bar{D} = \frac{\sum_{i=1}^n D_i}{n}$$

So, our first topic is paired t test. So, what does it do actually? If you want to describe it in simple language, so, sometimes we study the same group of subjects or elements in repeated samples and we may be interested to see how the mean is changing from one sample to the other. So here, as elements in one sample is matched with the element in another sample that is why it is called the paired sample or matched sample and we are going to now describe the t test, which could deal with these type of samples.

So, in many research designs, it is required to measure the same people more than once. So, here is an example, suppose manager of an office or your boss, he has decided to send some staffs for training so, that they can improve their typing speed. So, before he sent his staff to the training, the typing speed was recorded and then, after the training, when the person rejoins the job again the typing speed is recorded. So, if you send some 10, 15 or 20 staff for the training program, so, there is a before training mean typing speed and then there is an after training typing speed. So, you may be interested to see whether there is an improvement in the overall mean typing speed for this group of staff.

So, here, you have to apply paired t test and now I am going to briefly discuss about some features and assumptions. So, I have already mentioned these features that observations in the sample are appeared as they can be matched across the sample. So that we know we have discussed.

Now, let me state the assumptions clearly. So, we assume that population ideally be normally distributed, if it is not normally distributed, then at least large samples are required. And then, you assume that there are n pairs in the matched or related sample. So, it is better that you have n greater than or equal to 30.

So, to conduct the test, the first step is to compute the difference and this is for the i-th sample member. So here, I define this difference as d_i as x_{1i} minus x_{2i} . So x_1 and x_2 are basically 2 different groups. So, suppose this typing speed is a variable. So, 1 is basically after the training and 2 is before the training or you can think of several other examples,. But i is basically common, so, that is basically i-th sample member who is common in these 2 groups x_1 and x_2 .

So then, how do I define my point estimate for the population mean pair difference? So here is the formula. So, \bar{D} is the notation that I am introducing here as the point estimate and that is basically nothing you have to take the arithmetic mean of that variable d_i .

(Refer Slide Time: 04:21)

Paired t test

- ◆ Frame hypotheses:
 - ◆ H_0 : There is no significant difference between the averages of the two sets of sample $\rightarrow H_0: \mu_D = 0$
 - ◆ H_1 : There is significant difference between the averages of the two sets of sample $\rightarrow H_1: \mu_D < 0$ or, $H_1: \mu_D > 0$, or $H_1: \mu_D \neq 0$
- ◆ Define test statistic:
$$Z = \frac{\bar{D} - \mu_D}{\frac{\sigma_D}{\sqrt{n}}}$$
- ◆ If σ_D is unknown (most likely), estimate the unknown population std. dev. with a sample standard deviation:
$$S_D = \sqrt{\frac{\sum_{i=1}^n (D_i - \bar{D})^2}{n-1}}$$
- ◆ New statistic $t = \frac{\bar{D} - \mu_D}{\frac{S_D}{\sqrt{n}}}$ follows Student's t distribution with $n-1$ d.f.
- ◆ Follow the decision rules valid for t test

Now in the next step, you have to frame hypothesis as usual. So, here I start with the null hypothesis that there is no significant difference between the averages of that 2 sets of sample. And then, I say that okay, my H_0 can be represented as $\mu_D = 0$. So,

d is basically the difference variable that I have created and μ is basically the population mean.

Now, alternative hypothesis can be set in 3 different ways. So here, we define the alternative hypothesis as there is significant difference between the averages of the 2 sets of sample. And there can be 3 types of mathematical form that I can issue. So first I am going to give you a left-tail test, where you write alternative hypothesis as $\mu d < 0$. And the second type could be the right-tail test. And there you write your alternative hypothesis as $\mu d > 0$. Or you can have a two-tailed test where you can write H_1 as $\mu d \neq 0$.

So now in this case, you have to define the test statistic as the next step and here is the test statistic that is basically a z score. And then, you have \bar{d} that is the sample mean of the difference variable d . And then you have to take the difference between the sample mean of d and the population mean μd . So, if you work with the above shown null hypothesis then μd will of course take value 0, when you compute the test statistic.

And this difference now shall be divided by the σd divided by \sqrt{n} . So, the σd is of course, the unknown population standard deviation and most likely it will remain unknown. So, then what to do? So, you can actually get a proxy for σd . And the proxy could be defined as the sample standard deviation, which is defined as s_d . And I am showing you the formula in the red box. And that is not uncommon to you, you have seen how I can compute the sample variance. So, we are using the same formula, only in the place of x is now we are writing d .

So, once that s_d is computed, in the next step you write the expression for the new statistic and that is a t statistic. And I am showing you the expression in the box and this will follow a student's t distribution with $n - 1$ degrees of freedom. And then, you follow the decision rules valid for the t test that we have discussed before. It may also be interesting to go for a hypothesis testing for population proportion. So, sometimes it is interesting to compare the population proportions that come from 2 different populations.

So, you may get a sample from population 1 and you may get another sample from population 2. And you look at the sample proportion, they look different to you, but are they statistically different? That could be an interesting research question. So, in the next hypothesis testing, we are going to discuss this particular issue.

(Refer Slide Time: 07:50)

Testing Difference in Proportions

- Let two population proportions are π_1 and π_2
 $\bar{p} = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}$ (just average proportion)
 p_1 = proportion in group 1
 p_2 = proportion in group 2
 n_1 = number in group 1
 n_2 = number in group 2
- Let the point estimate for the difference is $p_1 - p_2$
- Test statistic is
$$Z = \frac{(p_1 - p_2) - (\pi_1 - \pi_2)}{\sqrt{\bar{p}(1-\bar{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

Recall, variance of a proportion is $p(1-p)/n$

←

<b style="color: red;">Lower-tail test: $H_0: \pi_1 - \pi_2 \geq 0$ $H_1: \pi_1 - \pi_2 < 0$	<b style="color: red;">Upper-tail test: $H_0: \pi_1 - \pi_2 \leq 0$ $H_1: \pi_1 - \pi_2 > 0$	<b style="color: red;">Two-tail test: $H_0: \pi_1 - \pi_2 = 0$ $H_1: \pi_1 - \pi_2 \neq 0$
--	--	--

- Conduct z test

So, we begin by introducing some notations. So, let 2 population proportions are π_1 and π_2 . And then, I define something as \bar{p} , that is a new thing that we are defining in this course and that is basically a weighted average of the sample proportions. So, here, if I am dealing with 2 groups, group 1 and group 2, the sample proportions are p_1 and p_2 . So, what is p_1 and what is p_2 ? So, if there are x_1 number of elements satisfying some criteria out of n_1 sample in group 1, then p_1 is defined as x_1 divided by n_1 . And similarly, I can define p_2 as well.

So, now, as there are n_1 and n_2 number of observations into different groups, respectively, you take weighted average and then you get the overall sample proportion \bar{p} . And then, once that is calculated, you now define your point estimate for the difference and now, you define that to be p_1 minus p_2 and next you decide on the test statistic and then you have the complicated z formula.

Well, it may look complicated initially, but it is not. Because if you recall certain things from the previous lectures, it is very common thing that we are still following here even in this particular slide. So, when you actually write a test statistic, you have to, in the numerator you have to first you talk about the difference between the sample statistic and the unknown population parameter value.

So, here p_1 minus p_2 is basically that point estimate that is basically as a statistic you can assume. And then, π_1 minus π_2 is basically the unknown difference between population proportions, π_1 and π_2 . So, the same philosophies or the concept is applied here in the numerator. And in the denominator, what are we seeing? So, if you remember, the variance of a proportion is basically p times 1 minus p divided by n , that we have seen earlier. So, now p is unknown. So, we have to get a proxy measure for p from the sample and that is basically your \bar{p} , the overall sample proportion.

So basically, now you have 2 different sample sizes n_1 and n_2 . So, you see the term that you are seeing under the square root that is basically the combined variance of 2 different samples. So of course, here also, we can see that there could be 3 different scenarios, 1 is lower-tail or left-tail test; 1 is upper or right-tail test; and there could be 1 two-tailed test.

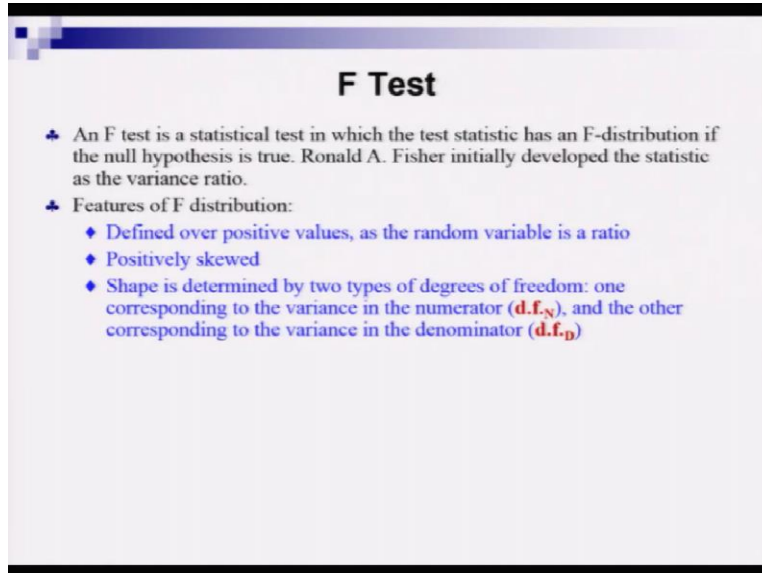
So, here first we are going to talk about the lower-tailed test. And I am showing you, how you can write the null and alternative hypothesis. So, nothing new here. Similarly, 1 can frame the null and alternative hypothesis for the upper-tail test, nothing new here as well. And then finally, you have that two-tailed test.

And you have to first decide which particular test you want to conduct. So, accordingly, you frame your hypothesis. And then, after that, you have to conduct a standard z test, because here, your test statistic is z . And as you are assuming that you have large samples so, n_1 and n_2 ideally is greater than or equal to 30. So, you can actually conduct as a z test.

So now, we are going to discuss the case of F test, we have not seen F test before. So, we are going to introduce this concept for the first time. And you will see later on in the course we will find very nice applications of d test. So, in the next slide, I am going to

show you the fundamental steps of conducting an F test. And I am not going to get into details of statistics regarding the F test, I am just going to talk about the most important things.

(Refer Slide Time: 12:20)



F Test

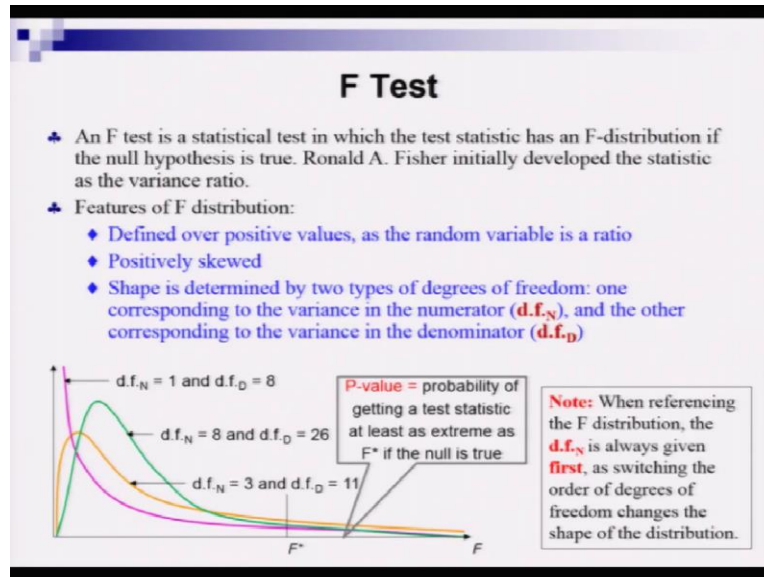
- An F test is a statistical test in which the test statistic has an F-distribution if the null hypothesis is true. Ronald A. Fisher initially developed the statistic as the variance ratio.
- Features of F distribution:
 - ♦ Defined over positive values, as the random variable is a ratio
 - ♦ Positively skewed
 - ♦ Shape is determined by two types of degrees of freedom: one corresponding to the variance in the numerator (**d.f._N**), and the other corresponding to the variance in the denominator (**d.f._D**)

So, let me first introduce the F test in the most simple possible language. So it is test for which the test statistic follows an F-distribution, if the null hypothesis is true and famous statistician, Ronald A. Fisher actually first developed the statistic as the variance ratio. If you remember, the way I defined the F value and the F distribution, while I was talking about the children of the normal distribution, then it was the ratio of the variances and that is, that was the form that Ronald Fisher introduced. And in his respect, this distribution finally was named F distribution.

So, there is another interesting feature that in most cases, the F distribution is a positively skewed distribution. But the shape of the F distribution is going to change because F is a very complicated distribution, it has 2 degrees of freedom.

Remember, in the chi-square case, we had 1 degrees of freedom. Now here, we have 2 types of degrees of freedom, because it is a ratio random variable. So, for the numerator, there is 1 degrees of freedom and for the denominator there is another degrees of freedom.

(Refer Slide Time: 13:39)



So, the shape is determined by d.f.N, which is the degrees of freedom corresponding to the variance in the numerator and the d.f.D, which is corresponding to the variance in the denominator of the ratio of 2 variances.

So, at the bottom part of the slide, I am going to show you a small diagram, a simple diagram where I am showing you various cases of F PDF. So, here you see the pink one which is looking almost like a rectangular hyperbola is F PDF for d.f.N equal to 1 and d.f.D equal to 8.

And you see if I now increase the d.f.N and d.f.D values to 3 and 11 respectively, then I get these orange colored PDF shape, which is a typical positively skewed distribution and if I continue to increase the d.f.N and d.f.D values to 8 and 26, then I get a PDF of higher height and it has got higher density in the right-tail also and that is given by the green colored PDF curve.

Now, I am going to tell you here about one statutory warning or you can see that as a cautionary note. So, when you are referencing F distribution, when you are trying to get the critical values from an F table, then d.f.N or the degrees of freedom corresponding to the numerator is always given first. So, you have to look at the column first and once you

fix the column for the numerator, degrees of freedom d.f.N, then you have to come to the row for d.f.D, which is basically the degrees of freedom for the denominator.

And this is very important. Why? Because switching the order of these d.f.N and d.f.D will have a big impact on the shape of the distribution and the corresponding critical values and probability values.

Now, in this context of F distribution, let us look at the very important concept of p value, once again. So here, now, let us concentrate on this diagram again. And now, we are going to talk about a specific value of F, which is the calculated value of test statistic. And let me denote by F star. So, now you erect a vertical line or draw a vertical line on F star. And then of course, it will partition your probability density function into 2 parts, 1 in the right hand side and 1 on the left hand side.

So now, when you concentrate on the area of right hand side, right with respect to F star, then basically that area below the F PDF curve, gives the p value. So, how do I interpret these p value? So here in this context, it is the probability of getting test statistic at least as extreme as F star if the null hypothesis is true.

(Refer Slide Time: 17:30)

Testing Difference in Variance

- A two-sample F test is used to compare two population variances when a sample is randomly selected from each population.
- Assumption: The populations must be independent and normally distributed.
- Types of F test:

Left tail Test
$H_0: \sigma_1^2 \geq \sigma_2^2$
$H_1: \sigma_1^2 < \sigma_2^2$

Right tail Test
$H_0: \sigma_1^2 \leq \sigma_2^2$
$H_1: \sigma_1^2 > \sigma_2^2$

Two tail Test
$H_0: \sigma_1^2 = \sigma_2^2$
$H_1: \sigma_1^2 \neq \sigma_2^2$
- S_1^2 and S_2^2 represent the sample variances with $S_1^2 > S_2^2$
- Degrees of freedom:
 - d.f._N = $n_1 - 1$ (corresponding to sample 1, associated with the higher sample variance)
 - d.f._D = $n_2 - 1$ (corresponding to sample 2, associated with the lower sample variance)

Now, we are going to look at 1 example of F test. And that is basically when we want to compare 2 population variances, we do not know the true values of the population

variances. But we want to draw some inference based on the sample variances. So how do you conduct a statistical testing to solve this problem?

So here are 2 sample F test is used to compare 2 population variances, when a sample is randomly drawn from each population. Here, we have to make a very strong assumption that the populations must be independent and normally distributed. And if you remember, our previous lectures, the F test could be of 3 types here, 1 of the 3 types, and they are namely left-tail test, right-tail test and two-tailed test. And I am showing you the null and alternative hypothesis under both categories, so you see, these expressions are given in these 3 boxes.

So, the left hand box talks about the left-tailed test and you see σ_1 , σ_2 , these are the notations I am using to denote the population standard deviations for population 1 and population 2. And if you follow that notation, then of course, these expressions are easy to follow. And then, in the center, I have right-tail test description of null and alternative hypothesis. And finally, at the right-hand side, the third box gives me the null and alternative hypothesis for a two-tailed test.

Now, the question is that we do not have any a priori knowledge about the population variance values. So, we do not know σ_1^2 and σ_2^2 . So then, how to proceed? We have to draw a sample and we have to calculate the sample variances from these 2 samples and let me call them s_1^2 and s_2^2 .

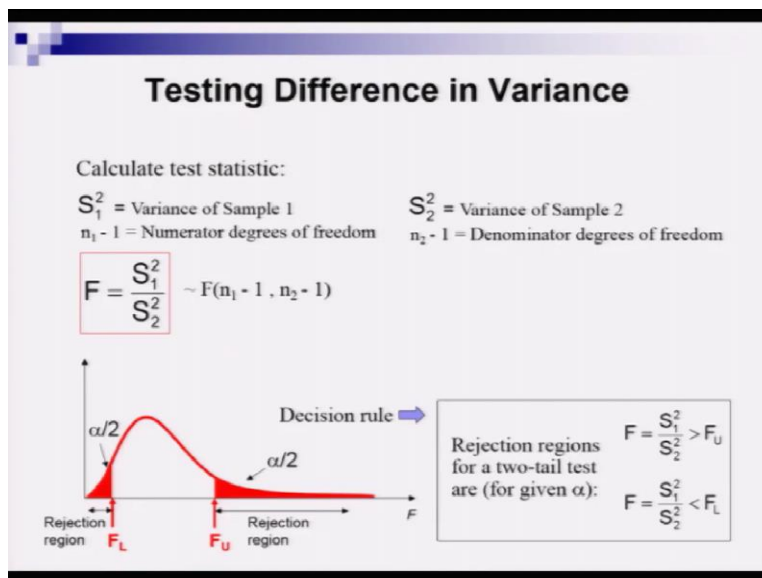
So, after calculating s_1^2 and s_2^2 , the sample variances most likely we are going to see that some, 1 of this is higher than the other. So, you call that higher variance corresponds to the population 1 so you name that population as population 1 for which you observe the highest value or the higher value of sample variance. And so, we can write $s_1^2 > s_2^2$. This is important to conduct an F test for checking the difference in variances.

Now, it is time to look at the degrees of freedom issue. So, if you remember that when we are conducting F test through F distribution, it has 2 degrees of freedom, 1 is associated with the numerator and other 1 is associated with the denominator. So, let us have a look

at them, what they are? So here d.f.N is n1 minus 1. So that is basically the number of observations in sample 1 minus 1 and d.f.D, that is basically corresponding to sample 2 and it is defined as number of observations in sample 2 minus 1.

So, I am not going to give you sentences where I am typing what to do next what to do next. So, here I am going to take a graphical approach to explain how to conduct an F test. We already have spoken about that diagram, 2 slides before. So, you will see a practical application of that slide again here in this slide, where I am going to describe this F test procedure again

(Refer Slide Time: 20:54)



So, this slide will start by calculating the test statistic value. So, we have to define the test statistic first. So, the test statistic is defined as capital F and that is the ratio of the sample variance from population 1 and sample variance of population 2. So, remember here that population 1 is that population, where the sample variance is higher when you compare 2 sample variances.

So, that is what we have written and these particular test statistic follows an F distribution with 2 degrees of freedom, namely in n1 minus 1 and n2 minus 1. So, I have actually explained them the previous slide. So next, you have to set the level of significance at alpha and then, you need to set the decision rule. So, for a particular level of alpha, as the

degrees of freedom are given for both numerator and denominator, it is not a difficult task to get the critical values, F_U and F_L . So, F_U is basically the critical value for the right-tail or it can be called critical value of upper-tail. And correspondingly we can find F_L that is basically the critical value for left-tail or the lower-tail and that is denoted as F_L .

So, now, I set a decision rule that if the calculated value of F is greater than F_U then you can reject your null hypothesis or if the calculated value of F is less than the critical value F_L , then also you can reject your null hypothesis. So, that is basically shown here in this diagram at the left-hand side and the bottom part of the slide.

So here, it is basically the same diagram that I have shown you, maybe a couple of minutes back. But here, I am indicating the rejection regions again, corresponding to or relative to F_L and F_U , the critical values and then, the task at hand is to get the F_L and F_U , in that case, you will be able to take a decision for this sample test.

Now, finding F values from statistical tables is a very challenging task, why? Because for a particular level of alpha, you have 2 degrees of freedom, so when you think about that matrix form of the F values, so in row you will have 1 set of degrees of freedom values, and on the columns you will have other set of degrees of freedom values. But note that, that matrix structure or the table will be valid for 1 particular alpha value. If you change alpha, then again, the same table will be reproduced but with different numbers in the sales as you have changed the alpha value.

So, F table it is very difficult to give you the entire picture in 1 slide, I am going to show you a glimpse of 1 particular F table for 1 particular level of significance. But there are many textbooks where you can find the longer versions of F table or you can use software or online calculators to find the critical values or the probability score in the F table. So, let me now take you to a short glimpse of a particular F table.

(Refer Slide Time: 24:48)

F Table for $\alpha = 0.05$

Rows show critical values for denominator degrees of freedom

Columns show critical values for numerator degrees of freedom

dof1 / dof2	2	4	6	8	10	12	15	20	30	40	60	120
2	19.000	19.247	19.330	19.371	19.396	19.413	19.429	19.446	19.462	19.471	19.479	19.487
4	6.944	6.589	6.463	6.411	6.384	6.372	6.358	6.346	6.335	6.328	6.323	6.320
6	5.143	4.824	4.724	4.683	4.663	4.654	4.647	4.642	4.638	4.636	4.635	4.634
8	4.459	4.178	4.097	4.066	4.051	4.044	4.039	4.036	4.034	4.033	4.033	4.033
10	4.103	3.857	3.792	3.767	3.755	3.750	3.746	3.744	3.743	3.743	3.743	3.743
12	3.885	3.665	3.614	3.594	3.586	3.583	3.581	3.580	3.580	3.580	3.580	3.580
14	3.739	3.543	3.498	3.481	3.476	3.474	3.473	3.473	3.473	3.473	3.473	3.473
16	3.634	3.461	3.421	3.406	3.402	3.400	3.399	3.399	3.399	3.399	3.399	3.399
18	3.555	3.400	3.364	3.351	3.348	3.347	3.346	3.346	3.346	3.346	3.346	3.346
20	3.493	3.353	3.321	3.309	3.307	3.306	3.306	3.306	3.306	3.306	3.306	3.306
22	3.443	3.317	3.288	3.277	3.276	3.275	3.275	3.275	3.275	3.275	3.275	3.275
24	3.403	3.291	3.264	3.255	3.254	3.254	3.254	3.254	3.254	3.254	3.254	3.254
26	3.369	3.263	3.238	3.230	3.229	3.229	3.229	3.229	3.229	3.229	3.229	3.229
28	3.340	3.240	3.216	3.209	3.208	3.208	3.208	3.208	3.208	3.208	3.208	3.208
30	3.316	3.220	3.197	3.191	3.190	3.190	3.190	3.190	3.190	3.190	3.190	3.190
40	3.232	3.146	3.124	3.118	3.117	3.117	3.117	3.117	3.117	3.117	3.117	3.117
60	3.150	3.075	3.054	3.049	3.048	3.048	3.048	3.048	3.048	3.048	3.048	3.048
120	3.072	3.007	3.000	3.000	3.000	3.000	3.000	3.000	3.000	3.000	3.000	3.000

Step 1: Find F_U from the F table for $n_1 - 1$ numerator and $n_2 - 1$ denominator degrees of freedom

Step 2: Find F_L using the formula: $F_L = 1/F_U^*$ Where F_U^* is from the F table with $n_2 - 1$ numerator and $n_1 - 1$ denominator degrees of freedom (i.e., switch the d.f. from F_U)

So here, I set alpha value add 0.05. And here, the rows show the critical values of the denominator, degrees of freedom. So, d.f.D and the columns show the critical values for the numerator of degrees of freedom. So, it is basically d.f.N. So, as I told you, 2, 3 slides before that you have to first look at the column and then, you have to come down rows and then you have to refer to the denominator degrees of freedom. So, that is the way you have to read the table.

But this is the table for right-tailed test. So, if you now have to conduct the two-tailed test, how do you proceed? it is actually not very simple. So, I am trying to express this in words. So now, look at the bottom of the slide, here I have written 2 steps. So, in step 2, you have to first identify a F_U from the F table for n_1 minus 1 numerator and n_2 minus 1 denominator degrees of freedom.

And then in step 2, you have to find the F_L , the critical value at the left-tail by using a formula and that is given by F_L equal to inverse of F_U^* . Now, this F_U^* is not the F_U that you have computed in step 1, this is somewhat different. So, what is it? So here, F_U^* actually is to be found from the F table with n_2 minus 1 numerator and n_1 minus 1 denominator degrees of freedom. So, you have to basically switch the degrees of freedom from the calculation of F_U .

So, we are going to end today's discussion, I thought to cover Cramér's V in today's lecture. But I had a second thought and I think that I will cover that concept in the lecture on the relationship between variables. So, wait for the next lecture. Thank you.