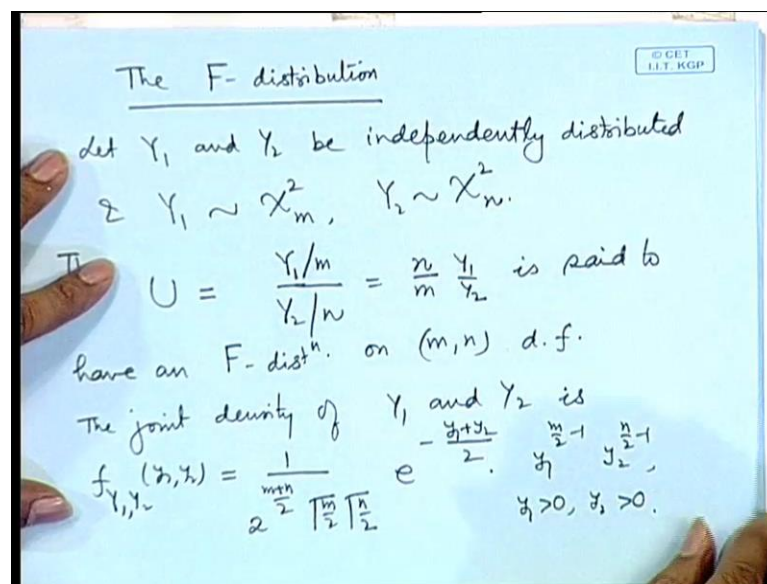


**Probability and Statistics**  
**Prof. Somesh Kumar**  
**Department of Mathematics**  
**Indian Institute of Technology, Kharagpur**

**Lecture - 48**  
**F-Distribution**

Sampling distribution which is important and it is used quite frequently is F- distribution.

(Refer Slide Time: 00:27)



So, now I introduce the F- distribution here. Let  $Y_1$  and  $Y_2$  be independently distributed random variables; and  $Y_1$  follow say chi square distribution on  $m$  degrees of freedom and  $Y_2$  follows say a chi square distribution on  $n$  degrees of freedom. Then if we define a variables called  $U$  as the ratio of this  $Y_1$  and  $Y_2$ , but divided by their degrees of freedom that is  $Y_1$  by  $m$  divided by  $Y_2$  by  $n$  that is basically becoming  $n$  by  $m$   $Y_1$  by  $Y_2$ , then this is said to have and F- distribution on  $m$   $n$  degrees of freedom. Now here one has to notice that two degrees of freedom terms are coming and therefore this order is important. If I am having a numerator chi square variable as  $m$  and the denominator chi square as  $n$  then we will write the ordered pair  $m$   $n$ ; that means, if I write  $n$   $m$ , it will denote a different F- distribution. Now by our theory of transformation of random variables,  $U$  is a function of  $Y_1, Y_2$ ; therefore, I can use a new dummy variable  $V$  and find out the joint density of  $U$  and  $V$  to derive the probability density function of  $U$ .

So, for that purpose we write the joint distribution of  $Y_1$  and  $Y_2$ ; the joint density of  $Y_1$  and  $Y_2$  is  $f$  of  $Y_1; Y_2$ . So, basically we multiply the individual densities of  $Y_1$  and  $Y_2$  which are basically chi square densities on  $m$  and  $n$  degrees of freedom. So, if we combined the coefficients  $1$  by  $2$  the power  $m$  plus  $n$  by  $2$ ,  $\Gamma$   $m$  by  $2$   $\Gamma$   $n$  by  $2$ ;  $e$  to the power minus  $y_1$  plus  $y_2$  by  $2$ ;  $y_1$  to the power  $m$  by  $2$  minus  $1$ ,  $Y_2$  to the power  $n$  by  $2$  minus  $1$ , where both  $y_1$  and  $y_2$  are positive.

So, we consider the transformation in which  $U$  is this variable.

(Refer Slide Time: 03:21)

Consider the transformation  
 $U = \frac{n}{m} \frac{Y_1}{Y_2}, \quad V = Y_2$   
 The inverse transformation is  
 $y_1 = \frac{m}{n} u v, \quad y_2 = v.$   

$$\begin{vmatrix} \frac{m}{n} v & \frac{m}{n} u \\ 0 & 1 \end{vmatrix} = \frac{m}{n} v$$
  
 So the joint density of  $U$  and  $V$  is  

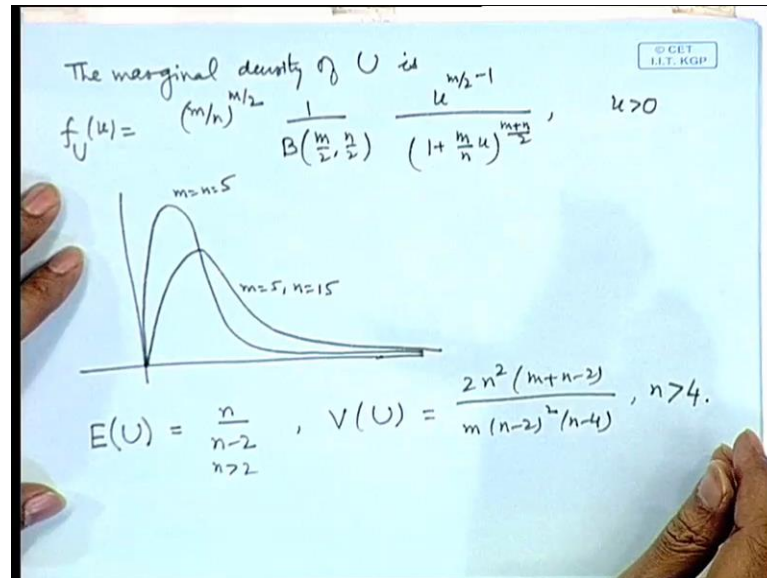
$$f_{U,V}(u,v) = \frac{\left(\frac{m}{n}\right)^{m/2}}{2^{\frac{m+n}{2}} \Gamma\left(\frac{m}{2}\right) \Gamma\left(\frac{n}{2}\right)} e^{-\frac{u}{2} \left(1 + \frac{m}{n} u\right)} \cdot u^{\frac{m}{2}-1} \cdot v^{\frac{m+n}{2}-1}$$
  
 $u > 0, v > 0$

So, consider the transformation  $U$  is equal to  $n y_1 / m Y_2$  and  $V$  is equal to say  $y_2$ . So, the inverse transformation is  $Y_1$  is equal to  $m$  by  $n$   $u$   $v$ ;  $Y_2$  is equal to  $v$ . So, the (Refer Time: 03:58) of the transformation can be calculated as  $m$  by  $n$   $v$   $m$  by  $n$   $v$   $u$ ;  $0$   $1$  which is basically  $m$  by  $n$   $v$  and since the terms are positive, so absolute value of  $j$  will also be the same. So, the joint density of  $U$  and  $V$  is; now you can observe, we are having here this constant term and we will be replacing  $y_1$  by  $m$  by  $n$   $u$   $v$  and  $y_2$  by  $v$ . So, here  $v$  by two terms will come out and this will give additional powers of  $v$ , for power of  $u$  will be this one only.

So after adjustment of the terms, we can write it as  $m$  by  $n$  to the power  $m$  by  $2$  divided by  $2$  to the power  $m$  plus  $n$  by  $2$   $\Gamma$   $m$  by  $2$ ,  $\Gamma$   $n$  by  $2$ ;  $e$  to the power minus  $v$  by  $2$ ;  $1$  plus  $m$  by  $n$   $u$ ;  $u$  to the power  $m$  by  $2$  minus  $1$ ;  $v$  to the power  $m$  plus  $n$  by  $2$  minus  $1$ , where  $u$  and  $v$  are positive variables. So, we can integrate with respect to  $v$  from

0 to infinity to get the desired density of  $f$  random variable. Again if we have observe the integral of this with respect to  $U$  is nothing, but a gamma function where the order is  $m$  by 2 and the coefficient is  $1$  plus half into  $1$  plus  $m$  by  $n$  u.

(Refer Slide Time: 05:51)

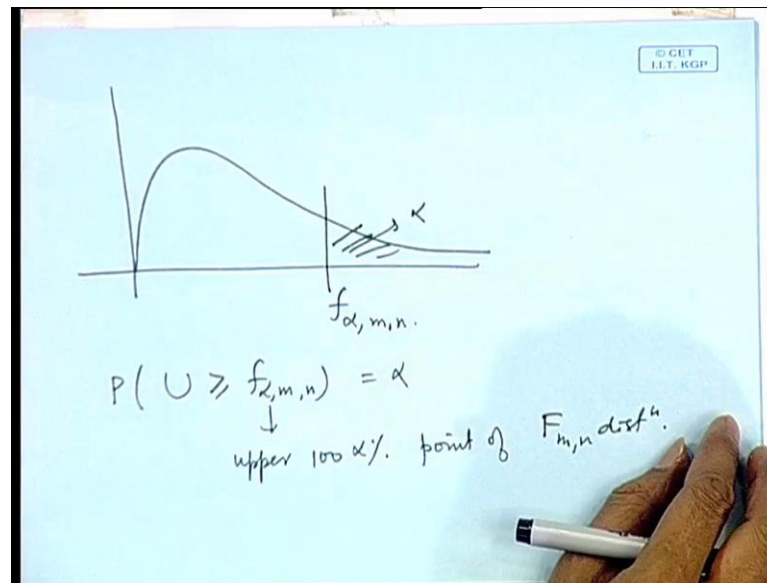


So using a standard argument, the marginal density will turn out to be the marginal density of  $U$  is then;  $f_U$  is equal to  $m$  by  $n$  to the power  $m$  by 2 and after the cancellation of the terms, we will get it as  $1$  by beta  $m$  by 2;  $n$  by 2;  $u$  to the power  $m$  by 2 minus 1 divided by  $1$  plus  $m$  by  $n$  u  $m$  plus  $n$  by 2 for  $u$  positive.

Obviously, this is a distribution of a positive value would random variable and it is positively skewed. However, the shape will vary depending upon the values of  $m$  and  $n$ . So, I can just give one example here; if I consider say  $m$  and  $n$  is equal to 5 then the form of the density is some what like this if I consider say  $m$  is equal to 5 and  $n$  is say 15, then the form is something like this. So, likewise for different values of  $m$  and  $n$  you get different shapes of the curves. Here calculation of the moments will make use of different beta functions. However, I will write the mean and variance, the mean of this is  $n$  by  $n$  minus 2; you may be little bit surprised here that it is dependent only upon second variable because  $n$  is not appearing here.

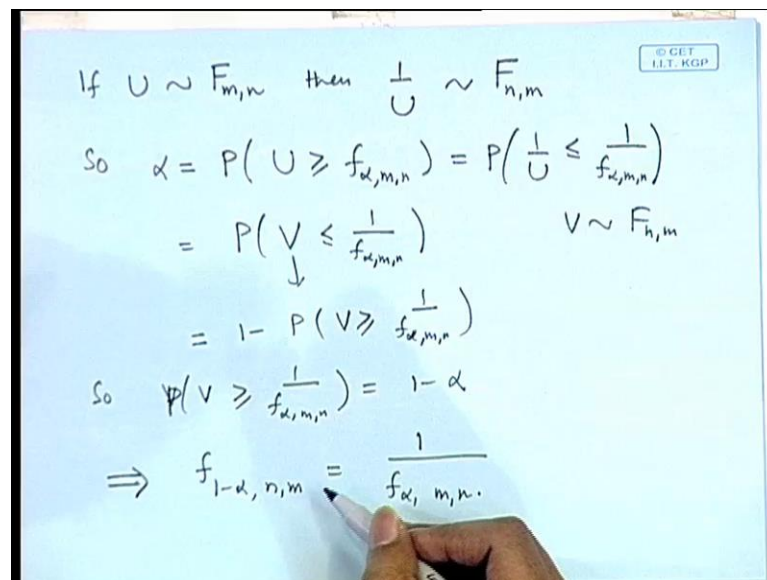
Variance of  $U$  is equal to here you need  $n$  greater then 2 and variance term is twice  $n$  square  $m$  plus  $n$  minus 2, divided by  $m$  into  $n$  minus 2 is square into  $n$  minus 4; this is valid for  $n$  greater then 4, this positively skew distribution.

(Refer Slide Time: 08:00)



So, we concentrate on the points here, so if I have this probability equal to alpha then this point is termed as  $f_{\alpha, m, n}$  that is the probability of  $U$  greater than or equal to  $f_{\alpha, m, n}$  is equal to alpha; that is this is upper 100 alpha percent point of  $F_{m, n}$  distributions.

(Refer Slide Time: 08:40)



Now by the definition of this  $f$  variable, it is clear that if I have  $U$  following  $F_{m, n}$  then  $1/U$  will follow  $F_{n, m}$  because  $1/U$  remains a ratio of chi square variables divided by their degrees of freedom; however, the numerator degrees of freedom have gone to the denominator and the denominator degrees of freedom have gone to the numerator.

So, this becomes F- distribution on n m degrees of freedom, so we can derive a formula for the points of F- distribution. So, we defined that  $f_{\alpha, m, n}$  is the upper 100 alpha percent point of the F- distribution on m n degrees of freedom. So, probability of U greater then or equal to this is equal to alpha. So, if I write it as probability of 1 by U less then or equal to 1 by  $f_{\alpha, m, n}$ ; then this 1 by U is a F n m variable that this is probability of sum V less then or equal to  $f_{\alpha, m, n}$ ; that is this v is F n m variable. So, if I am saying that the V greater then 1 by  $f_{\alpha, m, n}$ ; it is equality or inequality does not play any role here because of this continuous distributions. So, I am saying probability V greater then or equal to 1 by  $f_{\alpha, m, n}$  is equal to 1 minus alpha, but V is F n m distribution, this implies that  $f_{1-\alpha, n, m}$  is equal to 1 by  $f_{\alpha, m, n}$ .

This relationship is used for calculation of the percentage points of the F- distribution and generally the tables are because here it is a two dimensional table, you have m and n both varying and therefore, only for selected values of alpha; the tables are given. Now if they are given for say alpha is equal to 0.05 or alpha is equal to 0.1 then 1 minus alpha becomes 0.95 and 0.9 respectively. So, those tables can be automatically derived from the tables of 0.05 and 0.01 values etcetera.

(Refer Slide Time: 11:20)

The image shows a handwritten derivation on a blue notepad. It starts with two independent normal distributions:  $X_1, \dots, X_m \stackrel{i.i.d.}{\sim} N(\mu_1, \sigma_1^2)$  and  $Y_1, \dots, Y_n \stackrel{i.i.d.}{\sim} N(\mu_2, \sigma_2^2)$ . The sample variances are defined as  $S_x^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2$  and  $S_y^2 = \frac{1}{n-1} \sum_{j=1}^n (y_j - \bar{y})^2$ . Then, the scaled variances are shown to follow chi-square distributions:  $U = \frac{(m-1)S_x^2}{\sigma_1^2} \sim \chi_{m-1}^2$  and  $V = \frac{(n-1)S_y^2}{\sigma_2^2} \sim \chi_{n-1}^2$ . Finally, the ratio of these scaled variances is shown to follow an F-distribution:  $\frac{U/(m-1)}{V/(n-1)} = \frac{\sigma_2^2}{\sigma_1^2} \frac{S_x^2}{S_y^2} \sim F_{m-1, n-1}$ .

Now, we look at how in the sampling it arises, so if we consider say a random sample say  $x_1, x_2, \dots, x_m$  following normal  $\mu_1, \sigma_1^2$  and say  $Y_1, Y_2, \dots, Y_n$  be a

random sample form normal say  $\mu$   $\sigma^2$  is square and also I assume that these samples are taken independently.

Let us define the quantity say  $S_x^2$  as  $\frac{1}{m-1} \sum_{i=1}^m (X_i - \bar{X})^2$  and say  $S_y^2$  as  $\frac{1}{n-1} \sum_{j=1}^n (Y_j - \bar{Y})^2$ . Then by the theory of chi square  $\frac{S_x^2}{\sigma_1^2}$  follows chi square distribution on  $m-1$  degrees of freedom and  $\frac{S_y^2}{\sigma_2^2}$  follows, chi square distribution on  $n-1$  degrees of freedom; let me call its say  $U$  variable and this as the  $V$  variable. Then if I take the ratio  $U$  divided by  $m-1$  divided by  $V$  divided by  $n-1$ . Then this is nothing, but  $\frac{\sigma_2^2}{\sigma_1^2} \frac{S_x^2}{S_y^2}$  that follows  $F$ - distribution on  $m-1, n-1$  degrees of freedom.

So, this relationship or this result is used quite frequently in drawing inferences on ratios of the variances because ratios of the population variances and ratios of the sample variances is occurring here.

(Refer Slide Time: 13:23)

Thm: If  $T \sim t_n$ . Then  $T^2 \sim F_{1,n}$

$$T = \frac{X}{\sqrt{\frac{Y}{n}}}$$

$$T^2 = \frac{X^2}{Y/n}$$

$X^2 \sim \chi^2_1$   
 $Y \sim \chi^2_n$

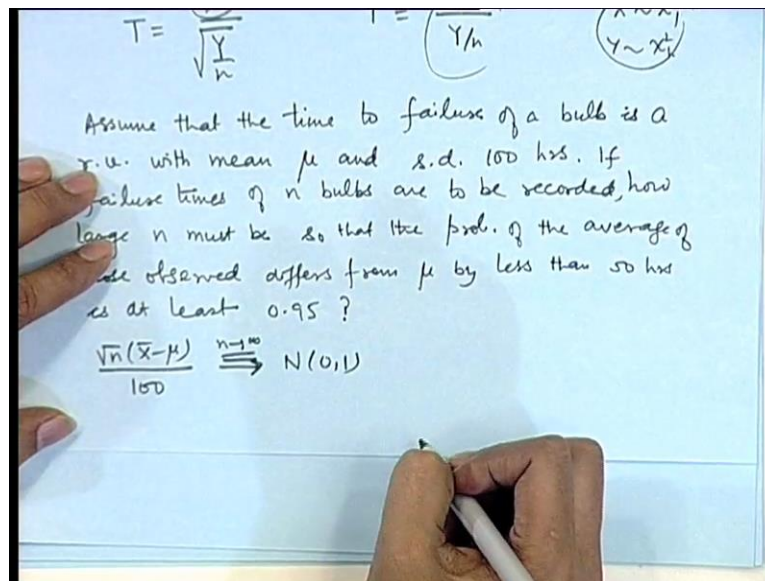
Another relationship which is coming here is that if say  $T$  follows  $t$  distribution on  $n$  degrees of freedom then  $T$  square follows  $F$ - distribution on 1 and  $n$  degrees of freedom. One can prove it by direct transformation by writing down the density of  $t$  and making the transformation  $U$  is equal to  $t$  square there and writes and compare with the forms of the densities. However, one can look at an easy representation see we can write  $T$  as  $X$



divided by a root  $\sqrt{Y}$  by  $n$ , where  $X$  is a standard normal and  $Y$  is a chi square variable, so if I look at  $t$  square that is  $X^2$  by  $Y$  by  $n$ .

Now, this  $X^2$  will be chi square on one degrees of freedom and  $Y$  will be chi square on  $n$  degrees of freedom. So, this is nothing, but the definition of an  $F$  variable on 1 and  $n$  degrees of freedom. Sampling distributions are extremely useful in a statistical inference and they are used all the time, I will give a couple of applications here.

(Refer Slide Time: 14:47)



Assume that the time to failure of a bulb is a random variable with say mean  $\mu$  and standard deviation say 100 hours. If failure times of  $n$  bulbs are to be recorded, how large  $n$  must be so that the probability of the average of those observed differs from  $\mu$  by less than 50 hours is at least 0.95.

That means how much should be my sample size such that the sample average and the population mean should differ  $\leq 50$  percent and this probability should be at least 0.95. So, we can make use of the central limit theorem here because the only information about the distribution that we are having is that; it is a particular distribution with certain mean and certain variance. So, the conditions for application of the central limit theorem are valid here. So, we will have a root  $n \times \bar{x} - \mu$  by  $\sigma$  that is 100; this will be approximately normal 0, 1 as  $N$  becomes large.

(Refer Slide Time: 16:58)

The image shows a whiteboard with handwritten mathematical steps. The steps are as follows:

$$P(|\bar{x} - \mu| \leq 50) \geq 0.95$$
$$\approx P\left(\frac{\sqrt{n}(\bar{x} - \mu)}{100} \leq \frac{50\sqrt{n}}{100}\right) \geq 0.95$$
$$\approx 2\Phi\left(\frac{\sqrt{n}}{2}\right) - 1 \geq 0.95 \quad \text{CLT}$$
$$\Rightarrow \Phi\left(\frac{\sqrt{n}}{2}\right) \geq 0.975$$
$$\Rightarrow \frac{\sqrt{n}}{2} \geq 1.96 \quad \text{or} \quad n \geq 16.$$

A small logo in the top right corner of the whiteboard reads "© CET I.I.T. KGP".

So, by this statement we have the condition here that probability of modulus  $\bar{x}$  minus  $\mu$  less than or equal to 50. We want to put this to be greater than or equal to 0.95, so we approximate this probability by converting to; by making use of the central limit theorem. So, this is less than or equal to  $50\sqrt{n}$  by 100, so this is approximately a standard normal random variable then this probabilities I can replace this variable by  $Z$ , where  $Z$  is a standard normal variable. So, this can be written in terms of the CDF that is twice  $\Phi$   $\sqrt{n}$  by 2 minus 1 greater than or equal to 0.95. So, I have made use of the central limit theorem here and this gives us  $\Phi$  of  $\sqrt{n}$  by 2 greater than or equal to 0.975. So, from the tables of the normal distribution  $\sqrt{n}$  by 2 must be greater than or equal to 1.96 or  $n$  must be greater than or equal to 16.

So, we need minimum sample size 16, so that the sample average and the population average do not differ by more than 50 and the probability of that should be at least 0.95. Let me give one more problem here.



(Refer Slide Time: 18:31)

The image shows a whiteboard with handwritten mathematical work. The work is as follows:

$$P\left(\frac{S_1^2}{S_2^2} > 3 \quad \text{or} \quad \frac{S_1^2}{S_2^2} < \frac{1}{3}\right) \quad \frac{S_1^2}{S_2^2} \sim F_{4,4}$$
$$= 1 - P\left(\frac{1}{3} \leq \frac{S_1^2}{S_2^2} \leq 3\right)$$
$$= 1 - \int_{1/3}^3 \frac{6x}{(1+x)^4} dx = 0.3125.$$

A hand holding a white marker is visible at the bottom of the whiteboard.

So, we consider independent random samples of size 5 from two normal populations and they have the same variance. So, what is the probability that the ratio of the larger to the smaller variance exceeds 3, so this problem is important in the following sense, see we have taken 2 samples from the same population. Now we want to check whether there is too much variability in the sampling process. So, we look at the variances of the 2 samples and one of them will be naturally larger and one will be smaller. So, we are saying that the ratio of the larger to the smaller exceeds 3 what is the probability of this event.

So, if we write in terms of  $S_1$  is square and  $S_2$  is square basically we are requiring here what is the probability that  $S_1$  square by  $S_2$  square is either greater than 3 or  $S_1$  square by  $S_2$  square is less than 1 by 3 what is the probability of this. Now if we have taken the samples of size 5 each, then by the formula that  $n - 1$  by  $m - 1$ ;  $S_1$  square by this 1, we get that  $S_1$  is square by  $S_2$  is square follows F- distribution on 4 and 4 degrees of freedom; that means, for calculation of this; we have to look at the either the tables of F 4 4 and we write down the density of F 4 4 here.

Fortunately the density of F 4 4 becomes a quite simple form. So, we can write this as 1 minus probability 1 by 3 less than or equal to  $S_1$  square by  $S_2$  is square less than or equal to 3 and this turns out to be 1 minus 1 by 3 to 3 and the density function of F 4 4 is 6 x

divided by 1 plus x to the power 4. So, this integration can be done easily and the value turns out to be 0.3125.

So, various problems which relate to the sample means or the sample variances or the comparison of the means or comparison of the variances can be solved using sampling distributions. In the portion of point estimation confidence interval estimation and testing of hypothesis, we will have frequent uses of these sampling distributions. So, in particular we have considered normal distribution itself has a sampling distribution because for the large sample, any sample mean will be approximately normally distributed under certain condition of course.

And then when we are sampling from normal distributions, then certain functions which are related to the means and the variances; they are having chi square, t, and F-distribution. So these are in particular 4 important sampling distributions. There are many more sampling distributions, but they are not as frequently used in practice.