

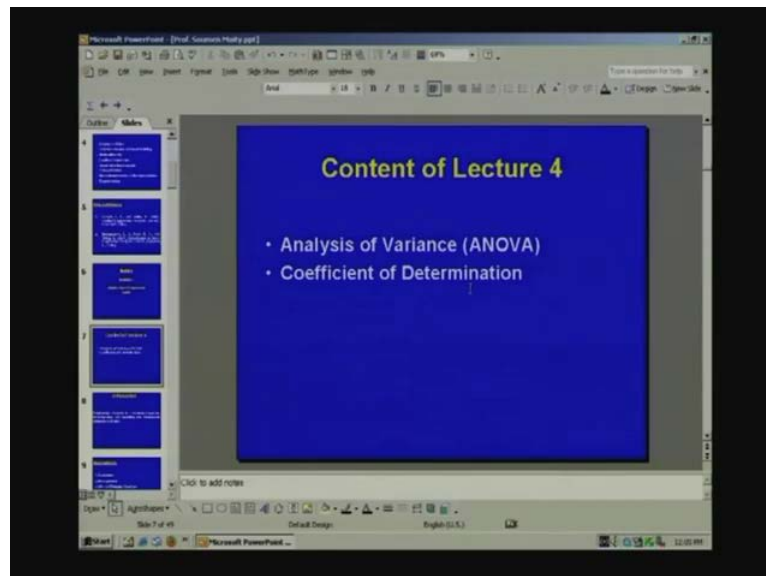
**Regression Analysis**  
**Prof. Soumen Maity**  
**Department of Mathematics**  
**Indian Institute of Technology, Kharagpur**

**Lecture - 4**  
**Simple Linear Regression (Contd.)**

Hi this is my 4th lecture on Simple Linear Regression, till now I mean given a set of observation, we know how to a beat Simple Linear Regression model to the data and once the module as been fitted then, we need to determine the goodness of the feet. So, and also we need to test, the statistical significance of the regression and coefficients.

Well so, one will do, this on like we in the last lecture, we have test it hypothesis is not, which is  $\beta_1$  equal to 0, against the alternative hypothesis that,  $\beta_1$  is not equal to 0. And we have used the test statistics  $t$  to test this hypothesis; so another way to approach this problem is the analysis of variance, so today we basically talking about ANOVA.

(Refer Slide Time: 02:28)



So, the content of lecture 4 is analysis of variance, in abbreviation it is ANOVA and also will talking about the coefficient of determination.

(Refer Slide Time: 02:53)

The whiteboard contains the following handwritten text:

$(x_i, Y_i)$   
 $i = 1(1)$

$Y = \beta_0 + \beta_1 X + \epsilon$   
 $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$

$H_0: \beta_1 = 0$  ag.  $H_1: \beta_1 \neq 0$

$t = \frac{\hat{\beta}_1}{\sqrt{\frac{MS_{Res}}{S_{xx}}}} \sim t_{n-2}$  under  $H_0$

Reject  $H_0$  if  
 $|t| > t_{\alpha/2, n-2}$

A small logo in the top right corner reads "© CEY I.I.T. KGP".

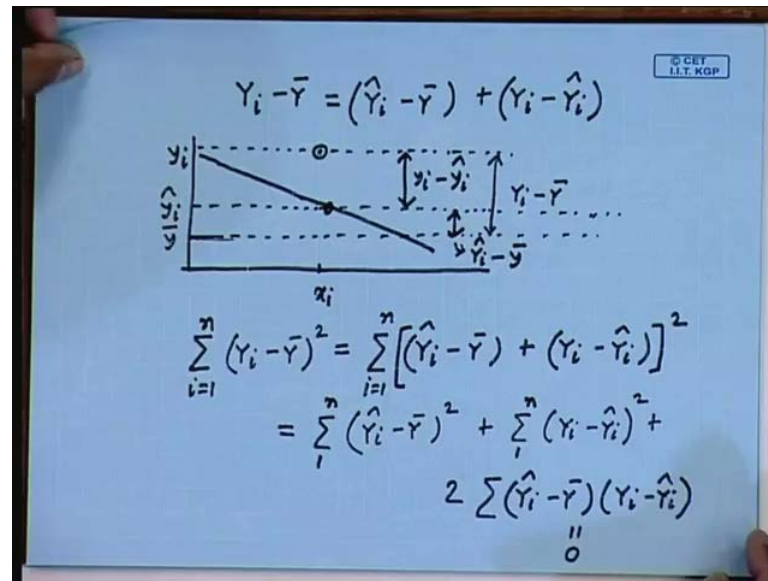
Well so, our model is  $Y$  equal to  $\beta_0$  plus  $\beta_1 X$  plus  $\epsilon$ , so given a set of data; say for example,  $X_i, Y_i$ , for  $i$  equal to 1 to  $n$ , we know how to fit a simple linear regression model to this data that is  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$ . And once the model has been constructed, it is important to confirm the goodness of fit and statistical significance of the regression coefficient.

So, we have tested the hypothesis  $H_0$ , which is equal to  $\beta_1 = 0$ , against the alternative hypothesis  $H_1$ , which says that  $\beta_1 \neq 0$ . And in lecture 3, we have used the test statistic  $t$ , which is equal to  $\hat{\beta}_1$  divided by the square root of  $MS_{Res}$  over  $S_{xx}$ , which follows  $t$  distribution with degree of freedom  $n-2$  and this is under  $H_0$ . So, we have used this test statistic to test the hypothesis  $H_0$  and we reject  $H_0$ , if the  $t$  value is greater than,  $t_{\alpha/2, n-2}$  at the level of significance, here is equal to  $\alpha$ .

So, now, another approach to solve this problem is called ANOVA technique, it is analysis of variance.



(Refer Slide Time: 10:49)



Now, this identity can be written as  $Y_i - \bar{Y}$ , which is equal to  $\hat{Y}_i - \bar{Y}$  plus  $Y_i - \hat{Y}_i$ . So, basically the significance of this identity is that, this is the deviation of the  $i$ th observation from overall mean. And this is the deviation of the  $i$ th observation from  $i$ th fitted observation, fitted value of  $i$ th observation from overall mean and this is the residual basically, this is  $e_i$ . So, how much of the variation, I mean, I am talking about the  $i$ th observation, so how much deviation of  $i$ th observation from overall mean explained by the model and this is the portion, which is remain unexplained.

Now, let me just draw figures, suppose given set of observation  $x_i, y_i$ , I have the fitted model, this is my fitted model and my  $i$ th observation is here, this is my  $i$ th observation. So, basically this is  $x_i$  and this height is  $y_i$ , now this is  $y_i$ , this is  $\hat{y}_i$ , because this point is basically  $x_i, \hat{y}_i$  and suppose, the overall mean of the responsible variable or of the data  $O_i$  is  $\bar{Y}$ , which is this well, now you see that, this is the distance is  $Y_i - \bar{Y}$ , so this is the deviation of  $i$ th observation from the overall mean.

Now part of the deviation is explained by the regression model and this distance basically, this distance is basically, it is  $\hat{y}_i - \bar{y}$  and this portion is  $y_i - \hat{y}_i$  well. So, the total deviation is this much and part of this deviation is explained by the regression model and remaining portion is the unexplained fault. Now, if we, if the square both side of this equation and sum from 1 to n, then we get summation  $Y_i - \bar{Y}$

$\bar{Y}$  whole square is equal to  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$  from  $i=1$  to  $n$ .

This is basically, the variation in the response variable or variance in the data, now you want split this variation into several parts, basically 2 parts the part, which is explained by the regression variable and the part the variance, which is not explained by the regression variable. So, you want to the part, which is not explained by the regression variable is basically  $SS_{residual}$  and we want to minimize the part, you want minimize is  $SS_{residual}$ .

You want the model to this as that, it can explain the variation in observation, I mean most of the part is explained by the model that is what, we want well. This is equal to  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$  plus 2 times  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})(Y_i - \hat{Y}_i)$  and I am going to prove that, this cross product on is equal to 0 well.

(Refer Slide Time: 17:52)

The image shows a handwritten derivation on a blue background. The derivation starts with the cross-product term (CPT) and shows it simplifies to zero. To the right, there are two equations defining  $\hat{Y}_i - \bar{Y}$  and  $Y_i - \hat{Y}_i$  in terms of the regression coefficients and variables.

$$\begin{aligned}
 \text{CPT} &= \sum_{i=1}^n (\hat{Y}_i - \bar{Y})(Y_i - \hat{Y}_i) \\
 &= \sum_{i=1}^n \hat{\beta}_1(x_i - \bar{x})[(Y_i - \bar{Y}) - \hat{\beta}_1(x_i - \bar{x})] \\
 &= \hat{\beta}_1 S_{xy} - \hat{\beta}_1^2 S_{xx} \\
 &= \hat{\beta}_1 (S_{xy} - \hat{\beta}_1 S_{xx}) \\
 &= 0
 \end{aligned}$$

$\hat{Y}_i - \bar{Y} = \hat{\beta}_1(x_i - \bar{x})$   
 $Y_i - \hat{Y}_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) = Y_i - \bar{Y} - \hat{\beta}_1(x_i - \bar{x})$   
 $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$

So, cross product on T is equal T summations  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})(Y_i - \hat{Y}_i)$ , now we can check that,  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})$  this is nothing but,  $\hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x})$  and similarly  $\sum_{i=1}^n (Y_i - \hat{Y}_i)$  is basically,  $\sum_{i=1}^n (Y_i - \bar{Y} - \hat{\beta}_1(x_i - \bar{x}))$ . So, just replacing  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})$  by the fitted value and now replace  $\sum_{i=1}^n (Y_i - \hat{Y}_i)$  that is basically  $\sum_{i=1}^n (Y_i - \bar{Y} + \hat{\beta}_1(x_i - \bar{x}))$ , so I can write this is equal to  $\sum_{i=1}^n (x_i - \bar{x})$ . So, if I

now replace 2 quantity here, i equal 1 to n, I will get some over beta 1 hat x i minus x bar into Y i minus Y bar.

So, here it is minus, minus beta 1 hat x i minus x bar, basically it is beta 1 hat S x y minus beta 1 hat square is S x x, this is the notation for this summation, which is equal to beta 1 hat S x y minus beta 1 hat is S x x. And this quantity is equal to 0, this is equal to 0, because we know that, beta 1 hat is equal to is x y by is S x x, so what we proves that the cross product on is equal to 0.

(Refer Slide Time: 21:56)

The image shows a handwritten derivation on a blue background. At the top, the total sum of squares is defined as  $\sum_1^n (Y_i - \bar{Y})^2 = \sum_1^n (\hat{Y}_i - \bar{Y})^2 + \sum_1^n (Y_i - \hat{Y}_i)^2$ . Below this, arrows point down to the standard notation:  $SS_T = SS_{Reg} + SS_{Res}$ . The regression sum of squares is further derived as  $SS_{Reg} = \sum_1^n (\hat{Y}_i - \bar{Y})^2 = \sum_1^n \hat{\beta}_1^2 (x_i - \bar{x})^2$ . A box highlights the final result:  $SS_{Reg} = \hat{\beta}_1^2 S_{xx}$ . Finally, the residual sum of squares is shown as  $SS_{Res} = \sum_1^n (Y_i - \hat{Y}_i)^2 = \sum_1^n e_i^2 \sim \chi_{n-2}^2$ .

So, we are left to it, then summation Y i minus Y bar whole square is equal to summation Y i hat minus Y bar whole square plus summation Yi minus Y i hat whole square. So, this quantity is denoted by S S T, so basically, it is a total sum of square and this quantity is denoted by S S regression; that means, sum of square due to regression and this is the portion, this called S S residual. So, what we have is that, S S total is equal to S S regression plus S S residual, while this is the splitting of total sum of square into 2 parts.

The total variation in Y is spitted into 2 parts, the variation the regression and the variation residual sum of square that means, variance which is not explained by the regression variable well. Now, we have been proved that, now what is S S residual S S sorry, S S regression is equal to summation Y i hat minus y bar whole square and this

quantity is nothing but, just now I mean, we have proved that, this is this quantity is equal to  $\beta_1 \hat{X}_i - \bar{X}$ .

So, square here square here, summation and this is going to be  $\beta_1^2 \sum x_i^2$ , so  $SS_{\text{regression}}$  is this quantity, now  $SS_{\text{residual}}$ , which is equal to summation  $(Y_i - \hat{Y}_i)^2$  from  $i=1$  to  $n$ , this is nothing but, the  $i$ th residual. So, I can write this own as summation  $e_i^2$  from  $i=1$  to  $n$  and you have proved that, this quantity is I mean, this follows chi square distribution with degree of freedom is not aim, it is the degree of freedom is  $n - 2$ , because you know all the  $e_i$  is there not independent.

We know that residual stay satisfy the constant that summation  $e_i$  is equal to 0 and the constant is summation  $e_i X_i$  is equal to 0. So, because of the due to this 2 constant all the  $e_i$  is not independent, you know out of  $n$   $e_i$  is you can choose  $n - 2$ ,  $e_i$  is independently and the remaining 2 have to be choose in such a way that, they satisfy this constant. So, be losing the 2 degree of freedom, because of these 2 constant, because  $e_i$  satisfies to constant that is why,  $SS_{\text{residual}}$  here, it follow chi square in minus anyway, I have proved this thing for also.

(Refer Slide Time: 27:21)

Handwritten mathematical derivation on a blue background:

$$SS_T = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

has DF  $(n-1)$

$$\begin{matrix} Y_1 - \bar{Y} \\ Y_2 - \bar{Y} \\ \vdots \\ Y_n - \bar{Y} \end{matrix}$$

$$\sum_{i=1}^n (Y_i - \bar{Y}) = 0$$

$$SS_T = SS_{\text{Reg}} + SS_{\text{Res}}$$

$$DF_T = DF_{\text{Reg}} + DF_{\text{Res}}$$

$$n-1 = 1 + (n-2)$$

Now,  $SS_T$  is sum of square, due to total sum of square, so  $SS_T$  is equal to summation  $(Y_i - \bar{Y})^2$  from  $i=1$  to  $n$  and this has the degree of freedom  $n - 1$ , because of the fact that, you know you have  $Y_1 - \bar{Y}$ ,  $Y_2 - \bar{Y}$  and



$\sum (Y_i - \bar{Y}) = 0$ ,  $i = 1$  to  $n$ . So, for this reason, you know out of  $n$  quantities  $n - 1$  can be chosen independently and the  $n$ th has to be chosen in such a way that, this constant is satisfied.

So, that is why, total sum of square has  $n - 1$  degree of freedom and also you know, you came to another  $SS_T = SS_{\text{regression}} + SS_{\text{residual}}$ , the variation, which as been explained by model. And this is the variation, which as not explained by the regression variable, if this as degree of freedom and also degree of freedom as the addit properties, so degree of freedom total is equal to degree of freedom of regression plus degree of freedom of residual.

So, this quantity is  $n - 1$  and we know that,  $SS_{\text{residual}}$  as degree of freedom  $n - 2$  and then, the degree of freedom of  $SS_{\text{regression}}$  is equal to 1.

(Refer Slide Time: 30:20)

**ANOVA TABLE**

Source of variation	DF	SS	MS	F	id.
Regression	1	$SS_{\text{Reg}}$	$MS_{\text{Reg}} = \frac{SS_{\text{Reg}}}{1}$	$F = \frac{MS_{\text{Reg}}}{MS_{\text{Res}}}$	
Residual	$n - 2$	$SS_{\text{Res}}$	$MS_{\text{Res}} = \frac{SS_{\text{Res}}}{n - 2}$		
Total	$n - 1$	$SS_T$			> ind

$E(MS_{\text{Reg}}) = \sigma^2$   
 $E(MS_{\text{Res}}) = \sigma^2 + \hat{\beta}_1^2 S_{xx}$

Now we make the ANOVA table, will source of variation degree of freedom sum of square well, let me write here, total sum of square and the source of variations are regression and the residual. So, this has degree of freedom 1, the residual has degree of freedom  $n - 2$  and this has degree of freedom  $n - 1$ , this one denoted by we know, what is this quantity, this denoted by  $SS_{\text{regression}}$ , this is called  $SS_{\text{residual}}$  and this is called  $SS_T$ .



Now M S mean square, which is obtained by dividing the S S by degree of freedom, so here it is M S regression is equal to S S regression by 1 and similarly, M S residual is equal to S S residual by degree of freedom that is n minus 2 yes. Now, we already know that expected value of M S regression is equal to sigma square, this we have proved before and it can be proved that, expected value of sorry, this is M S residual. It can be proved that M S regression is equal to sigma square plus beta 1 hat square S xx right.

(Refer Slide Time: 33:56)

$$\frac{(n-2)MS_{Res}}{\sigma^2} \sim \chi^2_{n-2}$$

$$\frac{MS_{Reg}}{\sigma^2} \sim \chi^2_1$$
 under  $H_0: \beta_0 = 0$

$$F = \frac{MS_{Reg}}{MS_{Res}} \sim F_{1, n-2}$$

To test  $H_0: \beta_0 = 0$   
 we compute  $F$  & Reject  $H_0$  if  $F > F_{\alpha, 1, n-2}$

Th. Let  
 $X \sim \chi^2_m$   
 $Y \sim \chi^2_n$

Then  

$$F = \frac{X/m}{Y/n} \sim F_{m, n}$$

And also, we know that n minus 2 M S residual by sigma square, this follows chi square with degree of freedom n minus 2 and it can be proved that, M S regression by sigma square, this follows, chi square with degree of freedom 1 that is under H naught, there is beta equal to 0. And these 2 quantities are this 2, these are basically this is function of random variables  $Y_i$ , so and they are independent.

Now, a statistical theorem, we have theorem says that, let X follows chi square m and Y follows chi square n and they are independent, then X by m by Y by n, this follows F with degree of freedom m and n, this quantity is divided by this. Basically, the residual of 2, chi square they follow, if distribution from there, from this theorem, we can say that, we can defined now, for 1 F for ANOVA table, it F is equal to M S regression by M S residual, this follows F with degree of freedom 1 n minus 2.

(Refer Slide Time: 36:54)

ANOVA TABLE

Source of Variation	DF	SS	MS	F	id.
Regression	1	$SS_{Reg}$	$MS_{Reg} = \frac{SS_{Reg}}{1}$	$F = \frac{MS_{Reg}}{MS_{Res}}$	
Residual	$n-2$	$SS_{Res}$	$MS_{Res} = \frac{SS_{Res}}{n-2}$		
Total	$n-1$	$SS_T$			> ind

$E(MS_{Reg}) = \sigma^2$   
 $E(MS_{Res}) = \sigma^2 + \hat{\beta}_1^2 S_{xx}$

And obviously, looking at this 2 expected value, so what we going to do is that ANOVA table next, we are going to compute, F this F is equal to M S regression by M S residual and looking at their expected value, it is initially clear that. If F is large, then it likely that, this beta 1 is not equal to 0, if beta 1 is equal to 0, then this residual going to be close to 1.

(Refer Slide Time: 38:04)

$\frac{(n-2) MS_{Res}}{\sigma^2} \sim \chi^2_{n-2}$   
 $\frac{MS_{Reg}}{\sigma^2} \sim \chi^2_1$

> ind.  
 under  $H_0: \beta_1 = 0$

$F = \frac{MS_{Reg}}{MS_{Res}} \sim F_{1, n-2}$

To test  $H_0: \beta_1 = 0$   
 we compute F & Reject  
 $H_0$  if  $F > F_{\alpha, 1, n-2}$

Pr. Let  
 $X \sim \chi^2_m$   
 $Y \sim \chi^2_n$

> ind

then  
 $F = \frac{X/m}{Y/n} \sim F_{m, n}$

So, from there, we can say that to test the hypothesis  $H_0$  beta 1 equal to 0, we compute F and reject not, if F is greater than, F alpha with degree of freedom 1 n minus

2. This is another way to another approach to test the hypothesis is beta 1 equal to 0, I mean this the same test, we can do using that, t distribution also basically, those 2 or same, I am going to proved that 1. First let me a give a 1 example for this ANOVA table well, I am going to constant the same example.

(Refer Slide Time: 39:37)

Handwritten notes on a blue background showing data, a regression equation, and ANOVA calculations.

Ad ( $x_i$ )     $Y_i$      $\hat{Y}_i = -0.1 + 0.7x_i$      $e_i$

1	1		
2	1		
3	2		
4	2		
5	4		

$SS_{Res} = 1.1$   
 $SS_T = \sum (Y_i - \bar{Y})^2 = 6$   
 $SS_{Reg} = \hat{\beta}_1^2 S_{xx} = (0.7)^2 \times 10 = 4.9$

Source of Variation	DF	SS	MS	F
Reg	1	4.9	4.9	13.6
Res	3	1.1	$\frac{1.1}{3} = 0.367$	
Total	4	6		

$F \sim F_{1,3}$   
 $F_{0.05, 1, 3} = 10.13$

Cost on advertisement that is  $X_i$  and this is the sales amount  $Y_i$  and we have the data 1 2 1 3 2 4 2 and 5 4 and we know the fitted model is  $\hat{Y}_i$  is equal to minus 0.1 plus 0.7  $X_i$ . And found here, we can you know before also, we compute it,  $e_i$  is and we know that,  $SS_{residual}$  for this problem is for this data is equal to 1.1 and what you need to compute is that, we need to compute. What is the total variation in the data  $SS_T$ , which is equal to summation  $Y_i$  minus  $\bar{Y}$  whole square, you can take that  $\bar{Y}$  is equal to 5 here and it is not difficult to check that, this is equal to 6.

And also we know that,  $SS_{regression}$ , which is equal to beta 1 hat squares  $S_{xx}$ , we know that beta 1 is 0.7. So, this is equal to 0.7 whole square and 1 can check that  $S_{xx}$ , basically we call to 10. So, this quantity is equal to 4.9 and here is my ANOVA table for this problem, ANOVA the source of variation, this is regression residual total degree of freedom, for this 1 is equal to 1 well, the total degree of freedom here is  $n - 1$  and  $n$  is equal to 5.

So, total degree of freedom is 4 and the degree of freedom residual is  $n - 2$  that is equal to 3 and hence the degree of freedom, for the regression is equal to 1 and the  $SS$

values are 4.9 1.1 and the total variation is equal to 6. So, here you can see that, for this problem the repeated model is really good, because the total variation in Y is 6 and most of the part of this variation has been explained by the regression.

So, out of 64.9 as been this is the part of the variation, which has been explained by the regression model, so most of the part has been explained the regression model and the portion, which is not explained by the regression model is 1.1 well. So, F value here is equal to 4.9 sorry, sorry, first step we need to compute M S value, the M S value is S S by degree of freedom. So, this is going to be 4.9 and M S residual 1.1 by 3, which is going to be 0.367 and the F value, is basically M S regression by M S residual, so 4.9 by 0.367, which is going to be 13.6.

And now you know, this F follows, this F here that follows, F distribution with degree of freedom 1 3. And you find the value of F 0.0513, which is going to be equal to 10 point 13, so you can take this value from the statistical table and now you check, you can see that our computed F value, which is 13.6 is larger than the tabulated value.

(Refer Slide Time: 45:48)

Handwritten calculations and ANOVA table:

$SS_{Res} = 1.1$   
 $SS_T = \sum (Y_i - \bar{Y})^2 = 6$   
 $SS_{Reg} = \hat{\beta}_1^2 S_{xx} = (0.7)^2 \times 10 = 4.9$

Source of Variation	DF	SS	MS	F
Reg	1	4.9	4.9	13.6
Res	3	1.1	$\frac{1.1}{3} = 0.367$	
Total	4	6		

$F \sim F_{1,3}$   
 $F > F_{0.05, 1, 3} = 10.13$   
 $H_0$  is Rejected

So, we can conclude that, which is I mean, this computed value is larger than F 0.051 with degree of freedom 1 and 3. So, we conclude that, we can reject H naught is rejected at the level of significance of this test is equal to alpha is equal to 0.05, so basically we have got the same result, using the t test. Now, I am going to prove that, this 2 test are

basically, you know this 2 whether you using F test or t test is does not mater in the case of simple linear regression, basically this 2 test are same well.

(Refer Slide Time: 47:04)

The whiteboard shows the following derivations:

$$F = t^2$$

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

$$t = \frac{\hat{\beta}_1}{\sqrt{\frac{MS_{Res}}{S_{xx}}}} = \frac{0.7}{\sqrt{\frac{0.367}{10}}} = 3.655$$

$$t_{n-2}^2 = \frac{\hat{\beta}_1^2 S_{xx}}{MS_{Res}} = \frac{MS_{Reg}}{MS_{Res}} = F_{1, n-2}$$

$$F = 13.61$$

$$t^2 = (3.655)^2 = 13.61$$

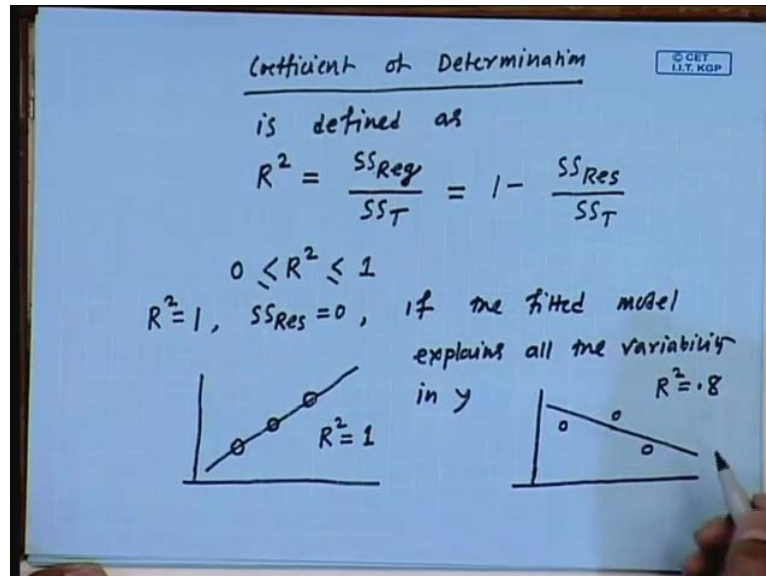
Let me prove that, F is nothing but, t square, I mean t square valuable well, so to test this hypothesis beta 1 is equal to 0, against the alternative the hypothesis beta 1 not equal to 0. Either you can go for t test, which is that t is equal to beta 1 hat square sorry, beta 1 hat by M S residual by S x x, this is done t statistic to test, this hypothesis, now we compute t square, which is going to be beta 1 hat square into S x x by M S residual right. And this quantity is well, this quantity is nothing but, M S regression and the trigonometry M S residual, which is nothing but, the F test.

So, the value of the distribution with degree of freedom 1 n minus 2 is same as the value of that distribution of the degree n minus 2, now just for you can check this 1, you know previous example, what we got is that for testing this hypothesis, we got F equal to 13.61. Now, if you go for the t statistic to test this hypothesis, we have this is beta 1 is equal to beta 1 add is equal to 0.07 and M S residual is 0.367 by S x x is equal t 10.

And this is going to be equal to 3.655 and you can check that, this 3, I mean that t square, which is basically equal to 3.655 is equal to 13.61, which is equal to F. So, whether you use the t statistic to test, this hypothesis or this is the F ANOVA approach that is the F statistic to test this hypothesis, there are basically same. And basically same for simply linear regression model once, we talking about multiple linear regression than, we need

to follow the ANOVA approach only. So, next I will be talking about the coefficient of determination well.

(Refer Slide Time: 51:02)



So, what is this that, this is denoted by  $R^2$  is defined as at square, which is basically ratio of  $SS_{Regression}$  by  $SS_{Total}$  well, you know several approach, to evaluate the performance of fitted model. So, this is 1 parameter, which can be, which is used to evaluate the performance of the repeated model, here this quantity is nothing but,  $1 - \frac{SS_{Residual}}{SS_T}$ .

And you know, that is  $SS_T$  is equal to  $SS_{Regression}$  plus  $SS_{Residual}$ , so that is why and obviously, the range for  $R^2$  is going to be from 0 to 1, it can be at most 1, if the  $SS_{Residual}$  equal to 0. So, this  $R^2$ , it is basically determines that, the proportion of variability that has been explained by the regression model, because  $R^2$  is equal to  $SS_{Regression}$  by  $SS_{Total}$ , so residual give you the proportion of variability that has been, you know explained by the regression model well.

Let me just give, this is the very important parameter, when  $R^2$  is going to be equal to 1, so  $R^2$  is going to be equal to 1, if  $SS_{Residual}$  equal to 0, so  $SS_{Residual}$  is equal to 0. And this will happened, if the predict model explained, this will happen, if the fitted model explains all the variability in  $Y$  right, that means, there is no parts, which remain unexplained by the model. And the example of this one is this is the case basically, suppose you have the data like this then, your fitted model is going to be this.

So here R square is going to be equal to 1 and this case, it is entirely clear that S S residual is going to be equal to 0, while if suppose in some example R square is equal to 0.8 may be, this is the case. These are the 3 data point and you fitted this model and entirely ability, you know is the R square is going to be very high.

And what we say here is that, what is the meaning of R square is equal to 0.8 is that, approximately 80 percent of the total variability in Y that has been explained by the regression model and then remaining 20 percent of the variability in Y remain unexplained by the model. So, the higher the value of the regress R square value, the better model is, so that is all for today and next class will be talking about, talking more about R square.

Thank you very much.