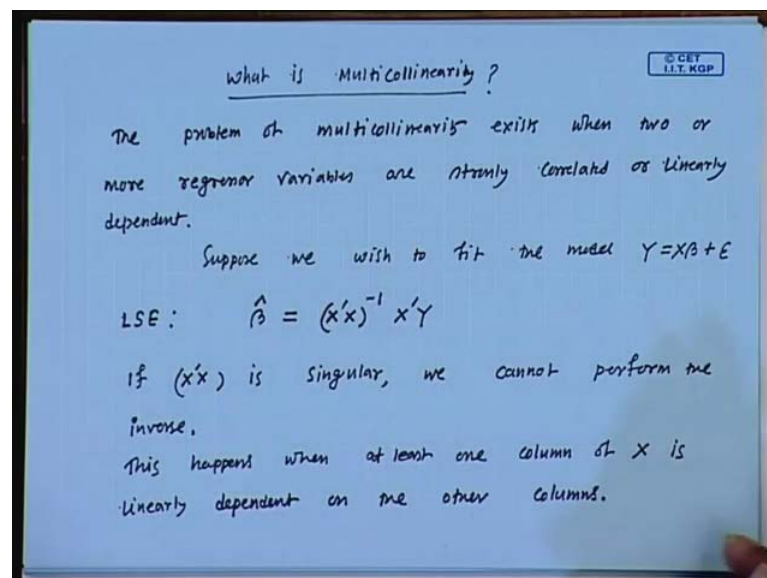


**Regression Analysis**  
**Prof. Soumen Maity**  
**Department of Mathematics**  
**Indian Institute of Technology, Kharagpur**

**Lecture - 14**  
**Multicollinearity**

This is my first lecture on Multicollinearity, and today we learn what is multicollinearity and the effects of multicollinearity or the problems due to multicollinearity.

(Refer Slide Time: 01:14)



So, what is multicollinearity? The problem of multicollinearity exists when two or more regressors variable are strongly correlated or linearly dependent. Suppose, we wish to fit the model  $Y$  equal to  $X$  beta plus epsilon, so we assuming that, we considering that, this is the matrix form of multiple linear regression and there are  $k$  minus 1 regressors and one response variable. So, we know that, the least square estimates of beta is beta hat equal to  $X$  prime  $X$  inverse  $X$  prime  $Y$ .

So now, if  $X$  prime  $X$  is singular then, we cannot perform the inverse,  $X$  prime  $X$  inverse and singular means, the determinant is equal to 0,  $X$  prime  $X$  determent determinant is equal to 0. And this happens, when atleast 1 column of  $X$  is linearly dependent on the other columns. So that means, the  $i$  th column of this matrix  $X$  stands for the  $i$  th regressor variable. So, the meaning of this one is that, if one regressor can be I mean, one

regressor is linearly dependent on the other regressors then,  $X'X$  is singular and we cannot compute the inverse of the  $X'X$  matrix, let me give one example for this one.

(Refer Slide Time: 05:58)

An Example

Can we use the data below to get a unique fit to the model  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$

$X_1$	$X_2$	$X_3$	$Y$
1	-2	4	81
2	-7	11	88
4	3	5	94
7	1	13	95
8	-1	17	123

$X'X$  is singular  
 $|X'X| = 0$

$\hat{\beta} = (X'X)^{-1} X'Y$

$X_2 + X_3 = 2X_1 \therefore X_1 = \frac{X_2 + X_3}{2}$

An example, can we use the data below to get a unique fit to the model  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$ . And the data carries data  $X_1, X_2, X_3$  and the response variable  $Y$ , so 1 -2 4 81 2 -7 11 88 4 3 5 94 7 1 13 95 8 -1 17 123. So, here in this example, we have three regressor variables  $X_1, X_2, X_3$  and one response variable, and here is the data. Now, the question is, whether we can use the data, use this data to get a unique fit to this model.

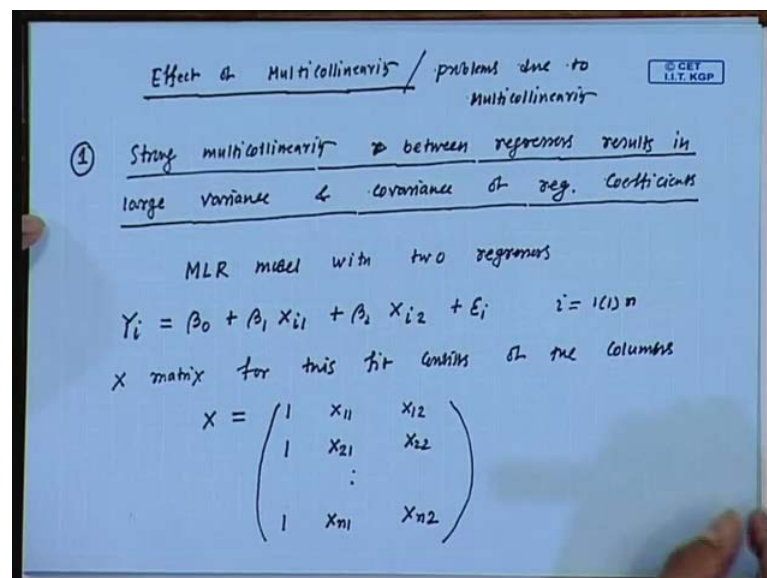
Of course, we can get a unique fit that means, we can estimate the regression coefficients using the least square estimate that is,  $\hat{\beta} = (X'X)^{-1} X'Y$ , provided  $X'X$  is not singular. Now here, if you observed, it is not difficult to check that, these three regressors are not independent. Here you check that,  $X_2 + X_3$ , take this first observation, so  $X_2$  value is -2,  $X_3$  value is 4 and  $X_1$  value is equal to 1, you can check that,  $X_2 + X_3$ , that is equal to 2, so  $X_2 + X_3$  is twice of  $X_1$ .

Take the second observation,  $X_2 + X_3$  is equal to 4, which is two times of  $X_1$ , third observation also  $X_2 + X_3$  equal to 8, which is  $2X_1$ . So here, the relation between the regressors is like  $X_2 + X_3$  is equal to twice  $X_1$  or same as writing that,  $X_1$  can be expressed in terms of using,  $X_1$  depends on  $X_2$  and  $X_3$ , so  $X_1$  is  $\frac{X_2 + X_3}{2}$ .

by 2. So here, really we cannot estimate the regression coefficients, because  $X'X$  is singular here, so here because of this fact, they are not independent.

So, that is why,  $X'X$  is singular, so  $X'X$  determinant is, here it is exactly equal to 0, determinant is equal to 0. So, this is an example here, which illustrates the definition of multicollinearity, here I mean, in this particular data, the problem of multicollinearity exists, because the columns of  $X$ , so the  $X$  is nothing but the  $X$  matrix is nothing but this one. So, here the one column for example, the first column can be I mean, it linearly depend can be written as the linear combination of  $X_2$  and  $X_3$ . So, that is why,  $X'X$  is singular and the problem of multicollinearity exists here. So, next we will be talking about the effect of multicollinearity.

(Refer Slide Time: 12:18)



Multicollinearity, the first or I can also say that, the problems due to multicollinearity, so the first one is, it says the strong multicollinearity between regressors, results in large variance and covariance of regression coefficients. So now, I am going to illustrate, what I mean by this, this says that, if there exist strong multicollinearity in the data then, that results in large variance and covariance of regression coefficients. So, let me consider multiple linear regression model with two regressors.

So, and my model here is  $Y_i$  equal to  $\beta_0$  plus  $\beta_1 X_{i1}$  plus  $\beta_2 X_{i2}$  plus  $\epsilon_i$ , so this  $i$  stands for the  $i$ th observation, so  $i$  is from 1 to  $n$ . Now, the  $X$  matrix for this fit consists of the columns, so here is the  $X$  matrix,  $X$  matrix for this model is, we

know that, this is  $1 \times 1 \times 1 \times 2 \times 1 \times 2 \times 1 \times 2 \times 2$  like that  $1, X_{n1} X_{n2}$ . Now, we will talk about the centering and scaling of this regression data.

(Refer Slide Time: 16:30)

Centering & Scaling Regression Data

if we write

$$x_{i1} = \frac{X_{i1} - \bar{X}_1}{\sqrt{S_{11}}}, \quad x_{i2} = \frac{X_{i2} - \bar{X}_2}{\sqrt{S_{22}}}, \quad y_i = \frac{Y_i - \bar{Y}}{\sqrt{S_{YY}}}$$

$$\bar{X}_1 = \frac{1}{n} \sum_{i=1}^n X_{i1}, \quad \bar{X}_2 = \frac{1}{n} \sum_{i=1}^n X_{i2}, \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$S_{11} = \sum_{i=1}^n (X_{i1} - \bar{X}_1)^2, \quad S_{22} = \sum_{i=1}^n (X_{i2} - \bar{X}_2)^2, \quad S_{YY} = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$X_1$	$X_2$	$Y$		$x_1$	$x_2$	$y$
$\bar{X}_1$				$\bar{x}_1 = 0$	$\bar{x}_2 = 0$	$\bar{y} = 0$
				$\sum_{i=1}^n x_{i1}^2 = 1$	$\sum_{i=1}^n x_{i2}^2 = 1$	$\sum_{i=1}^n y_i^2 = 1$
$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$				$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$		
$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}_1 - \hat{\beta}_2 \bar{X}_2$				$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}_1 - \hat{\beta}_2 \bar{x}_2$		

Centering and scaling of regression data, so what we do is that, if we replace, if we write  $X_{i1}$  is equal to  $X_{i1} - \bar{X}_1$  by  $\sqrt{S_{11}}$  square root of that. So, let me tell, what is this notation,  $\bar{X}_1$  of course, it is the mean associated with the first regressor. So, this is equal to  $\frac{1}{n} \sum_{i=1}^n X_{i1}$  and  $S_{11}$  is equal to  $\sum_{i=1}^n (X_{i1} - \bar{X}_1)^2$ . So, let me write also, say  $x_{i2}$  equal to  $\frac{X_{i2} - \bar{X}_2}{\sqrt{S_{22}}}$  and  $y_i$  is equal to  $\frac{Y_i - \bar{Y}}{\sqrt{S_{YY}}}$ .

So, this is the mean associated with the second regressor, so  $\bar{X}_2$  is equal to  $\frac{1}{n} \sum_{i=1}^n X_{i2}$  and  $S_{22}$  is equal to  $\sum_{i=1}^n (X_{i2} - \bar{X}_2)^2$ . And here also,  $\bar{Y}$  is equal to  $\frac{1}{n} \sum_{i=1}^n Y_i$  and  $S_{YY}$  is equal to  $\sum_{i=1}^n (Y_i - \bar{Y})^2$ . So, what I did here is that, my original data is  $X_1 X_2$  and  $Y$  and then, what I am doing is that, I am replacing this  $X_1$  by small  $x_1$ , so small  $x_1$  small  $x_2$ , just notation only and then, small  $y$ .

Now, one thing we need to observe is that, the mean of the original observation is  $\bar{X}_1$ . But here, if you take the mean of the transformed data, now this the

mean of small  $x_i$  is always equal to 0. So, here  $\bar{x}_1$  is equal to 0, similarly  $\bar{x}_2$  is also equal to 0 and  $\bar{y}$  is also equal to 0.

And the other thing to observe here is that, summation  $x_i^2$  from 1 to n, this is equal to 1 and this is also true for the second regressor like summation  $x_i^2$  from 1 to n, is equal to 1 and summation  $y_i^2$  from 1 to n, is equal to 1. Now, the model for the original data is, we wanted to fit this model  $Y$  equal to  $\beta_0$  plus  $\beta_1 X_1$  plus  $\beta_2 X_2$  plus epsilon, this transformation is called centering and scaling of regression data. And here, for this model we know that, the least square estimates gives  $\hat{\beta}_0$  is equal to  $\bar{Y}$  minus  $\hat{\beta}_1 \bar{X}_1$  minus  $\hat{\beta}_2 \bar{X}_2$ .

But, if you fit the same model here, for the transformed data also, if you fit the model like  $y$  equal to  $\beta_0$  plus  $\beta_1 x_1$  plus  $\beta_2 x_2$  plus epsilon. For the transformed data I am fitting the same model then, it is not difficult to check that, that  $\hat{\beta}_0$  is going to be 0 here, because see  $\hat{\beta}_0$  is equal to  $\bar{y}$  minus  $\hat{\beta}_1 \bar{x}_1$  minus  $\hat{\beta}_2 \bar{x}_2$ . But, all this you see,  $\bar{y}$  is equal to 0,  $\bar{x}_1$  is equal to 0,  $\bar{x}_2$  is also equal to 0. So here, this intercept for this transformed data is always going to be equal to 0, so that is why, for this scaled and centered data, we will omit the intercept  $\beta_0$  from the model.

(Refer Slide Time: 25:07)

The model, assuming that  $X_1, X_2 \perp Y$  are centered & scaled, is

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$$

$$X = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{pmatrix} = \begin{pmatrix} \frac{x_{11} - \bar{x}_1}{\sqrt{s_{11}}} & \frac{x_{12} - \bar{x}_2}{\sqrt{s_{22}}} \\ \frac{x_{21} - \bar{x}_1}{\sqrt{s_{11}}} & \frac{x_{22} - \bar{x}_2}{\sqrt{s_{22}}} \\ \vdots & \vdots \\ \frac{x_{n1} - \bar{x}_1}{\sqrt{s_{11}}} & \frac{x_{n2} - \bar{x}_2}{\sqrt{s_{22}}} \end{pmatrix}, Y = \begin{pmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix}$$

Normal equation

$$(X'X) \hat{\beta} = X'Y \Rightarrow \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} \sigma_{1y} \\ \sigma_{2y} \end{pmatrix}$$

So, the model assuming that,  $X_1 X_2$  and  $Y$  are centered and scaled is just  $y_i$  equal to  $\beta_1 x_i$  plus  $\beta_2 x_i$  plus epsilon  $i$ . So, we omitted the intercept  $\beta_0$  from

the model, because for the centered or the scaled data, we just checked that, intercept beta naught is always going to be equal to 0. So, next we will write down the X matrix for this model, so the X matrix is equal to, basically it is  $x_1$   $x_2$ , I am writing small that means, this is for the transformed data,  $x_1$   $x_2$  and  $x_n$   $x_n$ . See that, column 1 1 1 is not there, that column is associated with the intercept beta naught.

So, this one is basically,  $x_1$  is nothing but  $X_1$  minus  $\bar{X}_1$  by root of  $S_{11}$ ,  $x_2$  is capital  $X_2$  minus  $\bar{X}_2$  by root of  $S_{22}$ , and similarly  $x_n$  minus  $\bar{X}_n$  by root of  $S_{11}$   $X_n$  minus  $\bar{X}_n$  by root over of  $S_{22}$ . So, this is the X matrix for the transformed data you can say, now what is the normal equation associated with this model. The normal equation in general, basically I am trying to find the least square estimate of two regression coefficients, beta 1 and beta 2.

So, the normal equation is  $X'X\beta = X'Y$ , you know what is Y, Y is, in matrix form Y equal to capital Y 1 minus Y bar by root of  $S_{yy}$ , Y 2 minus Y bar by root over of  $S_{yy}$ . Similarly Y n minus Y bar by root over of  $S_{yy}$ , so this is what the y matrix is. So, here is the normal equation, now what is  $X'X$ , this one is, this implies say,  $X'X$  is, you can check that, the  $X'X$  is 1 and then,  $r_{12}$ , this  $r_{12}$  is nothing but the sample correlation between  $X_1$  and  $X_2$  and 1 here.

(Refer Slide Time: 31:07)

Centering & Scaling Regression Data

if we write

$$x_{i1} = \frac{X_{i1} - \bar{X}_1}{\sqrt{S_{11}}}, \quad x_{i2} = \frac{X_{i2} - \bar{X}_2}{\sqrt{S_{22}}}, \quad y_i = \frac{Y_i - \bar{Y}}{\sqrt{S_{YY}}}$$

$$\bar{X}_1 = \frac{1}{n} \sum_{i=1}^n X_{i1}, \quad \bar{X}_2 = \frac{1}{n} \sum_{i=1}^n X_{i2}, \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$S_{11} = \sum_{i=1}^n (X_{i1} - \bar{X}_1)^2, \quad S_{22} = \sum_{i=1}^n (X_{i2} - \bar{X}_2)^2, \quad S_{YY} = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$X_1$	$X_2$	$Y$	$x_1$	$x_2$	$y$
$\bar{X}_1$			$\bar{x}_1 = 0$	$\bar{x}_2 = 0$	$\bar{y} = 0$
			$\sum_{i=1}^n x_{i1}^2 = 1$	$\sum_{i=1}^n x_{i2}^2 = 1$	$\sum_{i=1}^n y_i^2 = 1$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}_1 - \hat{\beta}_2 \bar{X}_2$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}_1 - \hat{\beta}_2 \bar{x}_2$$

Because of the fact that, this sum  $x_i^2$  is equal to 1 and sum  $x_i^2$  is equal to 1, that is why these two elements are equal to 1.

(Refer Slide Time: 31:19)

The model, assuming that  $x_1, x_2 \leftarrow Y$  are centered & scaled, is

$$Y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$$

$$X = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{pmatrix} = \begin{pmatrix} \frac{x_{11} - \bar{x}_1}{\sqrt{s_{11}}} & \frac{x_{12} - \bar{x}_2}{\sqrt{s_{22}}} \\ \frac{x_{21} - \bar{x}_1}{\sqrt{s_{11}}} & \frac{x_{22} - \bar{x}_2}{\sqrt{s_{22}}} \\ \vdots & \vdots \\ \frac{x_{n1} - \bar{x}_1}{\sqrt{s_{11}}} & \frac{x_{n2} - \bar{x}_2}{\sqrt{s_{22}}} \end{pmatrix}, Y = \begin{pmatrix} \frac{y_1 - \bar{y}}{\sqrt{s_{yy}}} \\ \frac{y_2 - \bar{y}}{\sqrt{s_{yy}}} \\ \vdots \\ \frac{y_n - \bar{y}}{\sqrt{s_{yy}}} \end{pmatrix}$$

Normal equation

$$(X'X) \hat{\beta} = X'Y \Rightarrow \begin{pmatrix} 1 & r_{12} \\ r_{21} & 1 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} r_{1y} \\ r_{2y} \end{pmatrix}$$

And beta hat, here this matrix beta hat has two elements, beta 1 hat beta 2 hat, which is equal to  $r_{1y}$   $r_{2y}$ . This  $r_{1y}$  is the sample correlation between  $x_1$  and the response variable  $y$  and  $r_{2y}$  is the sample correlation between the second regressor and the response variable  $y$ .

(Refer Slide Time: 32:18)

Correlation matrix

$$\begin{pmatrix} 1 & r_{12} \\ r_{21} & 1 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} r_{1y} \\ r_{2y} \end{pmatrix}$$

where  $r_{12}$  is sample correlation between  $x_1$  &  $x_2$

$r_{1y}$  " " "  $x_1$  &  $y$

$r_{2y}$  " " "  $x_2$  &  $y$

$$r_{1y} = \frac{\sum_1^n (x_{i1} - \bar{x}_1)(y_i - \bar{y})}{\sqrt{s_{11} s_{yy}}}$$

$$r_{12} = \frac{\sum_1^n (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{\sqrt{s_{11} s_{22}}}$$

So, the normal equation we have is,  $r_{12}$  and  $r_{21}$  are same,  $\hat{\beta}_1$   $\hat{\beta}_2$ , is equal to  $r_{1y}$   $r_{2y}$ , where I already mentioned that,  $r_{12}$  or  $r_{21}$  is sample correlation between  $x_1$  and  $x_2$ , and  $r_{1y}$  is the sample correlation between  $x_1$  and  $y$ . Similarly,  $r_{2y}$  is the sample correlation between  $x_2$  and  $y$ , so what is this formula of  $r_{1y}$ ,  $r_{1y}$  is equal to  $\sum (x_{i1} - \bar{x}_1)(y_i - \bar{y}) / \sqrt{S_{11} S_{yy}}$  and  $r_{12}$  is equal to  $\sum (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) / \sqrt{S_{11} S_{22}}$ .

So, may be more precisely, I should replace them by capital X capital Y, so you can check with the original X matrix now. This is the X matrix for the transformed data and then, you can check, why  $X'X$  is equal to this, you know what is  $r_{12}$  and what is  $r_{1y}$   $r_{2y}$ , all these things. So, you can check that, why this element is  $r_{12}$  that is, the sample correlation between the regressor  $x_1$  and  $x_2$ . Now, to get  $\hat{\beta}_1$  and  $\hat{\beta}_2$ , we need to compute the inverse of this matrix, this matrix is called the correlation matrix now. This is the  $X'X$  matrix, this one is, now it is called the correlation matrix.

(Refer Slide Time: 36:35)

The handwritten derivation on the blueboard shows the following steps:

$$X'X = \begin{pmatrix} 1 & r_{12} \\ r_{21} & 1 \end{pmatrix}$$

Inverse of  $(X'X)$  is  $(X'X)^{-1} = \begin{pmatrix} \frac{1}{1-r_{12}^2} & \frac{-r_{12}}{1-r_{12}^2} \\ \frac{-r_{21}}{1-r_{12}^2} & \frac{1}{1-r_{12}^2} \end{pmatrix}$

$$\hat{\beta} = (X'X)^{-1} X'Y = \begin{pmatrix} \frac{1}{1-r_{12}^2} & \frac{-r_{12}}{1-r_{12}^2} \\ \frac{-r_{21}}{1-r_{12}^2} & \frac{1}{1-r_{12}^2} \end{pmatrix} \begin{pmatrix} r_{1y} \\ r_{2y} \end{pmatrix}$$

The estimates are

$$\hat{\beta}_1 = \frac{r_{1y} - r_{12} r_{2y}}{1-r_{12}^2} \quad \hat{\beta}_2 = \frac{r_{2y} - r_{21} r_{1y}}{1-r_{12}^2}$$

So, what we have is that, we have  $X'X$  is equal to  $\begin{pmatrix} 1 & r_{12} \\ r_{21} & 1 \end{pmatrix}$  here and you can check that, the inverse of  $X'X$  is  $(X'X)^{-1}$ , which is equal to  $\frac{1}{1-r_{12}^2} \begin{pmatrix} 1 & -r_{12} \\ -r_{21} & 1 \end{pmatrix}$ , it is not difficult to check that, this is the inverse of  $X'X$ .



So, beta hat, the least square estimator of the regression coefficient beta hat is equal to X prime X inverse X prime Y.

So, which is equal to  $\frac{1}{1 - r^2} \begin{bmatrix} 1 & r \\ r & 1 \end{bmatrix}^{-1} \begin{bmatrix} 1 & r \\ r & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$ , so from here, the estimates are beta 1 hat equal to  $\frac{y_1 - r y_2}{1 - r^2}$  and beta 2 hat is  $\frac{y_2 - r y_1}{1 - r^2}$ .

So, what I said the problem with, if there is a problem of multicollinearity in the data then, that results in large variance and covariance of correlation coefficients. So, right now, we have the least square estimator of the regression coefficient beta 1 and beta 2. Next we will check, what is the variance of beta 1 and we are going to prove that, the variance of beta 1 and also the variance of beta 2, that tends to infinity, as r tends to 1.

That means, when there is a strong multicollinearity between the regressor x 1 and x 2 then, the variance of beta 1 hat and the variance of beta 2 hat are going to be infinity, so that is what, we are going to prove now. So, next we need to find the variance of beta 1 hat and variance of beta 2 hat in terms of the sample correlation coefficient r.

(Refer Slide Time: 41:14)

$$V(\hat{\beta}_j) = \sigma^2 (X'X)^{-1}_{jj}$$

$$V(\hat{\beta}_1) = \sigma^2 \cdot (X'X)^{-1}_{11} = \frac{\sigma^2}{1 - r_{12}^2} \rightarrow \infty \text{ as } |r_{12}| \rightarrow 1$$

$$V(\hat{\beta}_2) = \sigma^2 (X'X)^{-1}_{22} = \frac{\sigma^2}{1 - r_{12}^2} \rightarrow \infty \text{ as } |r_{12}| \rightarrow 1.$$
 If there is strong multicollinearity between  $x_1$  &  $x_2$ , then the correlation coefficient  $r_{12}$  will be large.
 
$$\text{Cov}(\hat{\beta}_1, \hat{\beta}_2) = \sigma^2 (X'X)^{-1}_{12} = \frac{-\sigma^2 r_{12}}{1 - r_{12}^2} \rightarrow \pm \infty$$
 depending on whether  $r_{12} \rightarrow +1$  or  $r_{12} \rightarrow -1$ .

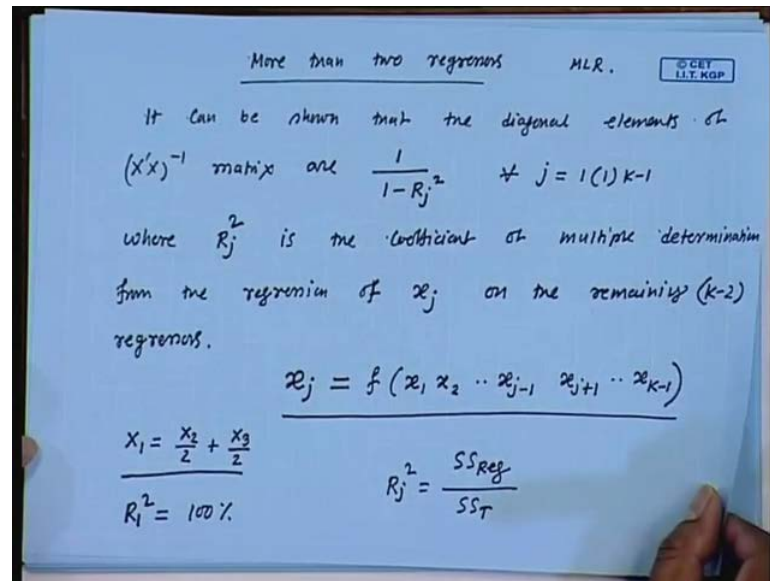
So, we know that, variance of  $\hat{\beta}_j$  in general, is equal to  $\sigma^2 X'X^{-1}$   $j$ th element of this  $X'X^{-1}$ . So here, from this formula, variance of  $\hat{\beta}_1$  is equal to  $\sigma^2$  into  $X'X^{-1}$  1 1 th element. So, what is  $X'X^{-1}$  1 1 th element, the 1 1 th element is  $1/(1 - r_{12}^2)$ , so this one is equal to  $\sigma^2/(1 - r_{12}^2)$ . Now, so this is the variance of  $\hat{\beta}_1$  and similarly, the variance of  $\hat{\beta}_2$  is equal to  $\sigma^2 X'X^{-1}$  2 2 th element, which is also equal to, 2 2 th element is also  $1/(1 - r_{12}^2)$ .

So here, the variance of  $\hat{\beta}_2$  is also  $\sigma^2/(1 - r_{12}^2)$ , now if there is strong multicollinearity between the regressor  $x_1$  and  $x_2$  then, the correlation coefficient  $r_{12}$  will be large. So, large means, the modulus value can be equal to 1, so when the modulus of  $r_{12}$  tends to 1, the variance of  $\hat{\beta}_1$  tends to infinity, so this tends to infinity, as  $r_{12}$  tends to 1. So, this will tend to 1, when the regressor  $x_1$  and  $x_2$  are strongly correlated or you can say, when there is strong multicollinearity between  $x_1$  and  $x_2$  then,  $r_{12}$  tends to 1.

Similarly, this quantity is also, the variance of  $\hat{\beta}_2$  also tends to infinity, as  $r_{12}$  tends to 1. So, also I said that, strong multicollinearity results in large variance and covariance, so what is a covariance of  $\hat{\beta}_1$  and  $\hat{\beta}_2$ , this is equal to  $\sigma^2 X'X^{-1}$  1 2 th element. So, what is 1 2 th element or 2 1 th element that is,  $-r_{12}/(1 - r_{12}^2)$ , so this one is going to be  $\sigma^2 r_{12}/(1 - r_{12}^2)$  square.

And this also tends to plus minus infinity depending on, whether  $r_{12}$  tends to plus 1 or  $r_{12}$  tends to minus 1. So, this is the proof of I mean, we using two regressors in the model, in the multiple linear equation model, we illustrated how the strong multicollinearity results in large variance and covariance of the regressor coefficients, so the same. So, this illustration is using the two regressors in the multiple linear regression model, but this is also true, if you have more than two regressors in the model, so that we are going to mention now.

(Refer Slide Time: 47:15)



Suppose, you have more than two regressors in the multiple linear regression model then, it can be proved that, it can be shown that, the diagonal elements of  $X'X$  inverse matrix are  $1/(1 - R_j^2)$ , for all  $j$  from 1 to  $k - 1$ . So, I am talking about multiple linear regression model with  $k - 1$  regressors and of course, I need to define, what is this  $R_j^2$ . So here, before we had  $R_{12}^2$ , now  $R_{12}^2$  has been replaced by  $R_j^2$ , where  $R_j^2$  is the coefficient of multiple determination from the regression of  $x_j$  on the remaining  $k - 2$  regressors.

I think I need to explain this, here  $R_j^2$  is the coefficient of multiple determination, we know what is coefficient of multiple determination in multiple linear regression model that is,  $R^2$  is equal to  $SS_{Reg} / SS_T$ . And this parameter measures the proportion of variability in the response variable, that is explained by the regressor variable. But here, what I mean by  $R_j^2$  is that, so that is, we know what is  $R^2$ , the coefficient of multiple determination, when we fit a regression model between  $y$  and the, between the response variable and the  $k - 1$  regressor variables. But here, this  $R_j^2$  is the coefficient of multiple determination from the regression of  $x_j$  on the remaining  $k - 1$  regressors. That means, here the multiple linear regression model is in between  $x_j$  and remaining  $k - 1$  regressors. So here,  $x_j$  is expressed in terms of the remaining regressors say,  $x_1, x_2, \dots, x_{j-1}, x_{j+1}, \dots, x_{k-1}$ .

So, here the model is inbetween, trying to express  $x_j$  in terms of the remaining regressors. So, if you can recall my first example today, there we have checked that, the  $X_1$ , if I am recalling correctly,  $X_1$  is equal to  $X_2$  by 2 plus  $X_3$  by 2 and this one is true for all the observations we had, so this is what we want. And the coefficient of multiple determination for this example, this one is basically  $R^2$  and here,  $R^2$  is 100 percent, because  $X_2$  and  $X_3$  can explain 100 percent of the variability in  $X_1$ .

So, this  $R^2$ , now I hope you understood that what is  $R^2$ ,  $R^2$  is the coefficient of multiple determination for this regression model. So, you fit the model, you express  $x_j$  in terms of the remaining regressors and then, you compute  $SSE$  regression, for this model you compute  $SSE$  total and you will get  $R^2$ .

(Refer Slide Time: 53:50)

Handwritten notes on a whiteboard:

$$V(\hat{\beta}_j) = \sigma^2 (X'X)^{-1}_{jj}$$

$$= \frac{\sigma^2}{1-R_j^2} \rightarrow \infty \text{ as } R_j^2 \rightarrow 1$$

If there is strong multicollinearity between  $x_j$  & any subset of other  $(k-2)$  reg, then  $R_j^2$  will be close to unity.

So, since this is the diagonal element, now I can say that, when you have more than two regressors in the multiple linear regression model, the variance of  $\hat{\beta}_j$  is equal to  $\sigma^2 (X'X)^{-1}_{jj}$  and the  $(j,j)$  element is equal to  $\sigma^2 / (1 - R_j^2)$ . So, this tends to infinity, as  $R_j^2$  tends to 1,  $R_j^2$  tends to 1 I mean, sometime we write, most of the time we write this in terms of percentage, so  $R_j^2$  is 100 percent means that,  $R_j^2$  tends to 1.

So,  $R_j^2$  tends to 1 means, this will happen, if there is strong multicollinearity between  $x_j$  and any subset of other  $k - 2$  regressors then,  $R_j^2$  will be close to unity. So, if there is a strong multicollinearity between the regressor  $x_j$  and any subset of

the other remaining  $k - 2$  regressors then,  $R^2_j$  will be close to unity and the variance of  $\hat{\beta}_j$  will tend to infinity.

So, this proves, just today we could manage to talk about only one effect of multicollinearity that means, only one problem due to multicollinearity. That says that, strong multicollinearity results in large variance and covariance of the regression coefficients. And we illustrated this point, both for the multiple linear regression model using two regressor variable in the model and then, in general. Next class, we will be talking about some other, I know there are several effects of multicollinearity in multiple linear regression model. So, we will be talking about those problems in the next class.

Thank you very much.