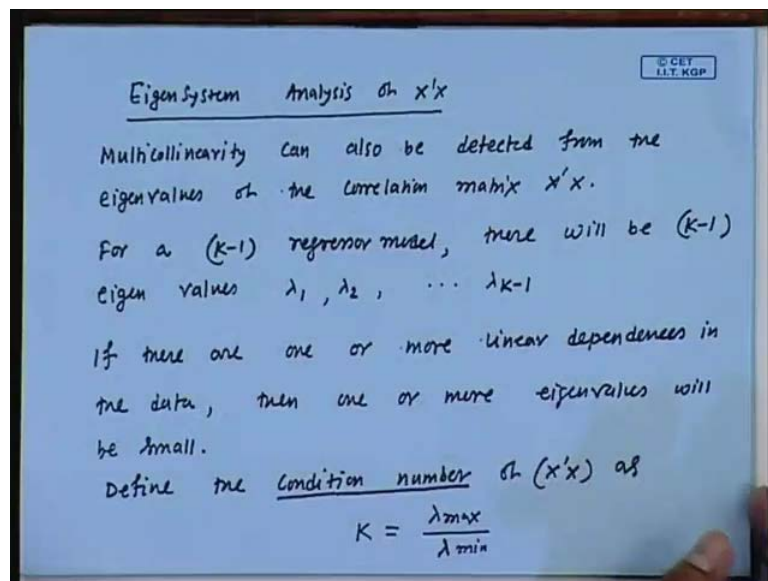


**Regression Analysis**  
**Prof. Soumen Maity**  
**Department of Mathematics**  
**Indian Institute of Technology, Kharagpur**

**Lecture - 16**  
**Multicollinearity (Contd.)**

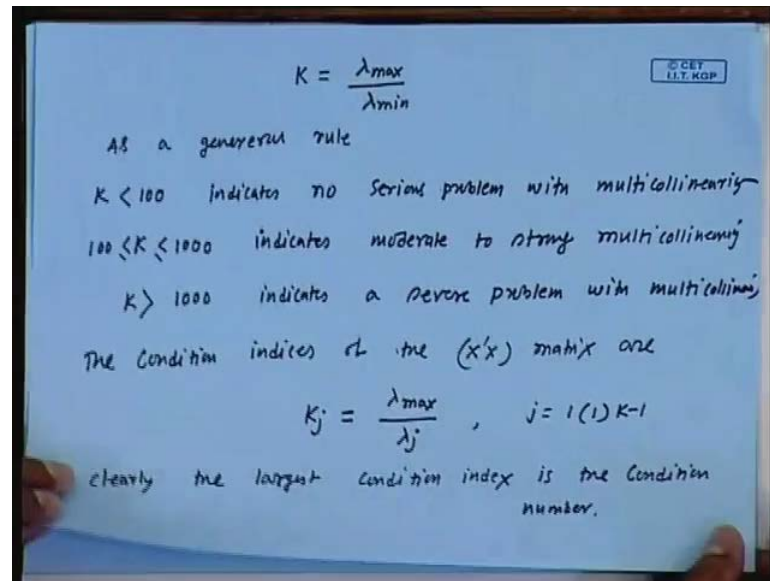
Hi, this is my 3rd lecture on Multicollinearity. The problem of multicollinearity exists when two or more regression variables are dependent or you can say that when two or more random regressive variables are linearly dependent. So, in the last class, we talked about different technique techniques to detect multicollinearity, we learned about examination of correlation matrix, and also we learned about you know Eigen system analysis of  $X'$  prime index matrix or the correlation matrix.

(Refer Slide Time: 01:59)



So, first I will recall the Eigen system analysis we learned in the previous class, well what we will do here is that we first compute the  $K$  minus 1, Eigen values of the  $X'$  prime  $X$  matrix, and then we compute the condition number which is  $K$  denoted by  $K$  equal to  $\lambda_{\max}$  by  $\lambda_{\min}$ .

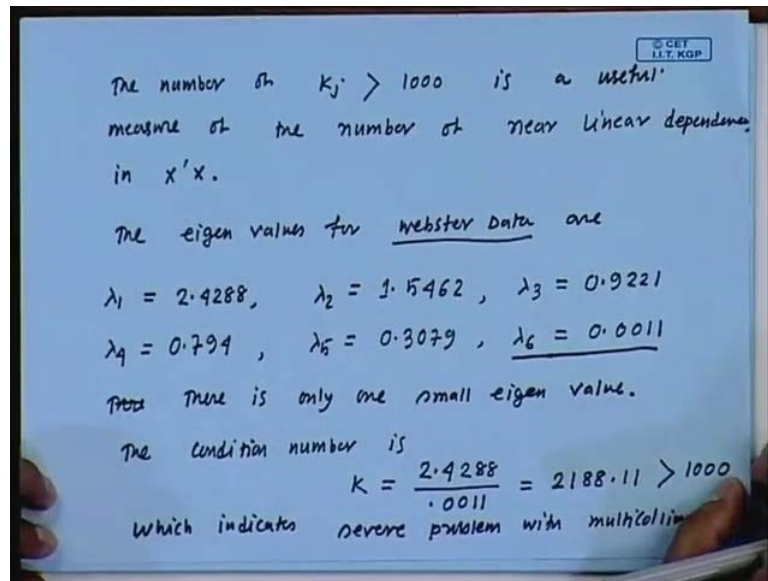
(Refer Slide Time: 02:36)



And the larger value of  $K$  indicates the severe problem with multicollinearity as a general rule you know this is what we talked in the last class, if  $K$  is less than 100 then that indicates no serious problem with multicollinearity. If  $K$  is between 100 and 1000 then that indicates moderate to strong multicollinearity, but if  $K$  is greater than 1000 then that indicates a severe problem with multicollinearity.

So, this is how we you know condition number is a used to detect the multicollinearity, the advantage of this Eigen system analysis is that it not only detect multicollinearity, it can measure the number of linear dependencies in the correlation matrix. And also, you know it can determine or it can identify the nature of linear dependencies between the regresses, so for that you know what we do is that we compute the condition indices.

(Refer Slide Time: 04:40)



Here, is the condition indices  $K_j$ ,  $K_j$  is associated with the  $j$ th regressor, so the  $K_j$  is lambda maximum by lambda  $j$ , and the number of  $K_j$  greater than 1000 is a useful measure of the number of linear dependencies in  $x$  prime  $x$ .

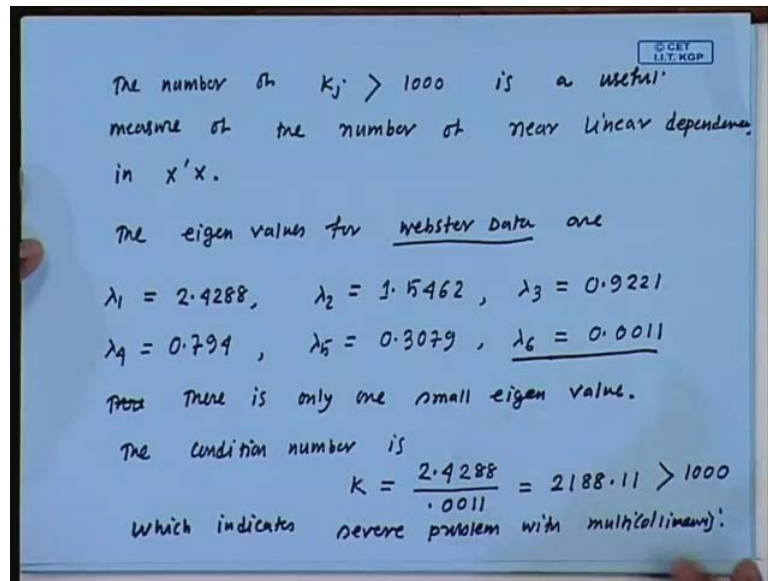
(Refer Slide Time: 05:05)

Unstandardized Regressor & Response Variables from Webster, Gunst, and Mason (1979)

Y	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>
10.006	8	1	1	1	0.541	-0.099
9.737	8	1	1	0	0.130	0.070
15.057	8	1	1	0	2.116	0.115
8.422	0	0	9	1	-2.397	0.252
8.625	0	0	9	1	-0.046	0.017
16.289	0	0	9	1	0.365	1.509
5.958	2	7	0	1	1.996	-0.865
9.313	2	7	0	1	0.228	-0.055
12.960	2	7	0	1	1.930	0.502
5.541	0	0	0	10	-0.798	-0.399
8.756	0	0	0	10	0.257	0.101
10.937	0	0	0	10	0.440	0.432

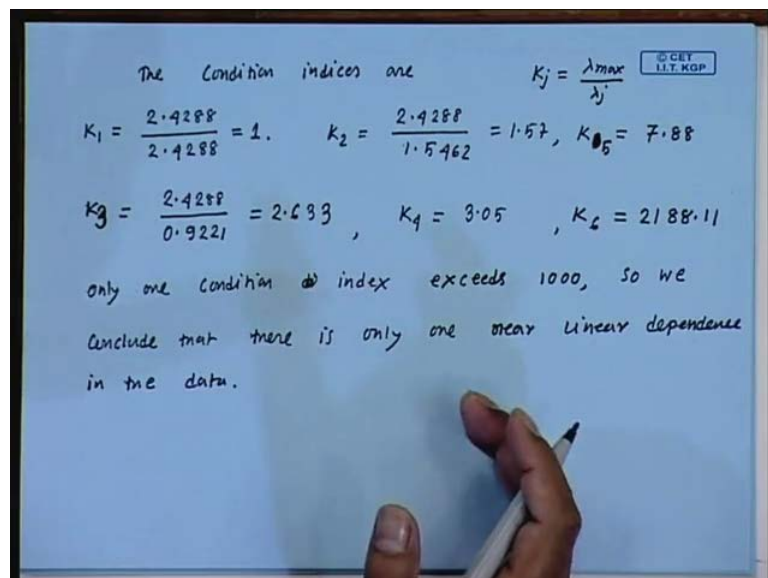
And then I illustrated this result using the webster data, this is the webster data we talked about in the previous class, we have 6 regressors response variable right and I referred this data as you know webster data.

(Refer Slide Time: 05:29)



Now, for the Eigen values of  $X'X$  matrix of the correlation matrix for the Webster data is here, you know  $\lambda_1 \lambda_2 \lambda_3 \lambda_4 \lambda_5 \lambda_6$  and the smallest Eigen value is this one, which is closed to 0. Now, we compute the condition number, so the condition number is 2188 which is greater than 1000, so this condition number indicates the presence of severe multicollinearity in the Webster data. So, next what we will do is that, we will compute the condition indices also since we know the Eigen values we can compute the condition indices.

(Refer Slide Time: 06:40)

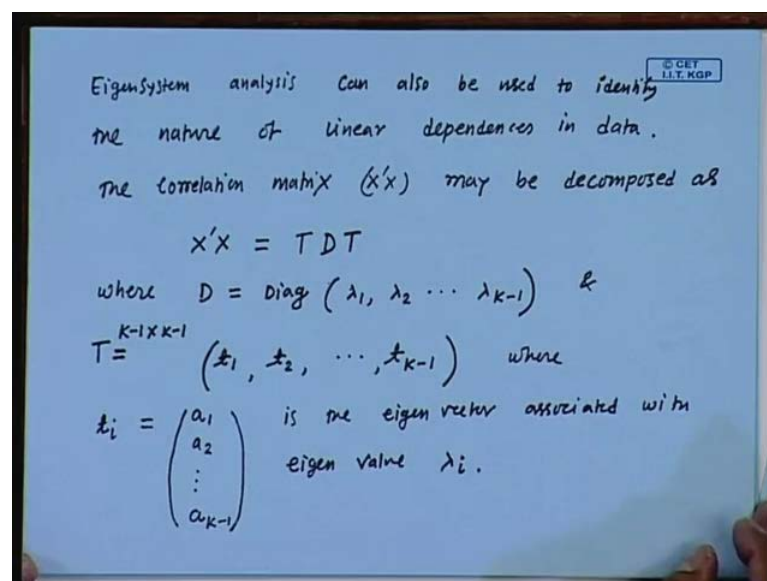


So, the condition indices are what is condition indices  $K_j$ ,  $K_j$  is equal to  $\lambda_{\max} / \lambda_j$ , so  $K_1$  is  $\lambda_{\max} / \lambda_1$  is basically maximum, so  $\lambda_1$  is 2.428 which is equal to 1. The  $K_2$  is  $\lambda_{\max} / \lambda_2$  which is equal to 2.4288 by  $\lambda_2$  which is equal to 1.5462 equal to 1.57. Similarly, you compute  $K_3$  which is equal to  $\lambda_{\max} / \lambda_3$  is equal to 7.88,  $K_4$  is  $\lambda_{\max} / \lambda_4$  that is 2.4288 by  $\lambda_4$  is 0.9221 which is equal to 2.633.

So, I did this a mistake here, this is  $K_3$  this is  $K_5$ , any way you can compute for this  $K_3$  and then  $K_4$  is equal to 3.05,  $K_5$  is equal to this quantity and  $K_6$  is equal to 2188.11. Now, we know that here only one condition index exceed 1000, that is  $K_6$ , so we conclude that there is only one near linear dependents in the data, because you know we I mentioned before that the number of  $K_j$  greater than 1000 that measures the number of linear dependencies in the data.

And since, here only one  $K_j$  that is  $K_6$  is greater than 1000, that is why the number of linear dependents or dependencies is equal to one only. Well, now what will this technique has lot of advantages like you know it not only detect the presence of multicollinearity it can measure the number of linear dependencies in the data and also it can identify the nature of the linear dependencies.

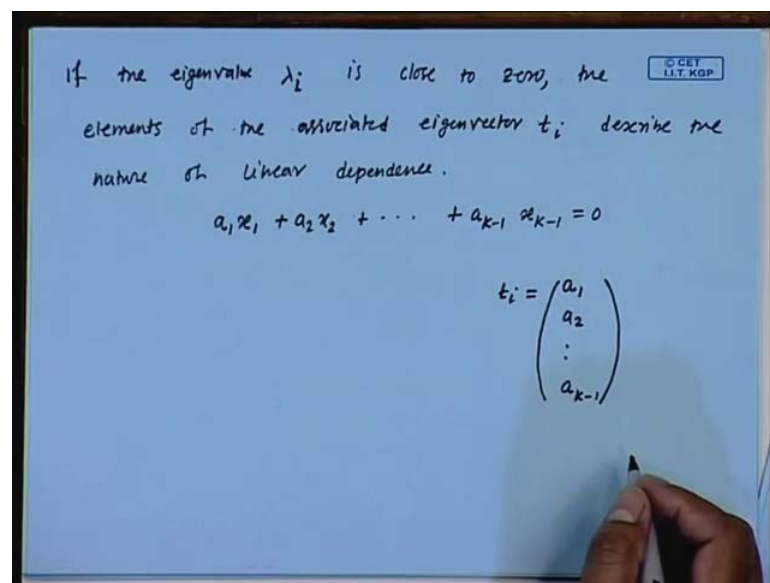
(Refer Slide Time: 11:03)



So, now, we will explain that portion, here Eigen system analysis can also be used to identify the nature of linear dependencies in data. So, let me explain this portion first, the

correlation matrix  $X'X$  may be decomposed as you know  $X'X$  you can write or decompose as  $X'X$  is equal to  $TDT$ , where  $D$  is a diagonal matrix, whose main diagonal elements are the Eigen values. So,  $D$  equal to diagonal  $\lambda_1 \lambda_2 \dots \lambda_{K-1}$ , and  $t$  is equal to  $t$  is a  $(K-1) \times (K-1)$  matrix, and whose the columns of this  $T$  matrix their  $t_1 t_2 \dots t_{K-1}$ . So, here this  $t_i$  is the Eigen vector associated with  $\lambda_i$ , where  $t_i$  is equal to  $a_1 a_2 \dots a_{K-1}$  is the Eigen vector associated with Eigen value  $\lambda_i$ , so this is the decomposition of the correlation matrix.

(Refer Slide Time: 15:23)



Now, if  $\lambda_j$ , if the Eigen value  $\lambda_i$  is closed to 0, the elements of the associated Eigen vector, that is  $t_i$  describe the nature of linear dependents. So, the nature of this linear dependents is like  $1 \times 1 + 2 \times 2 + \dots + (K-1) \times (K-1)$  is equal to 0, so this coefficient of the regressor variables  $a_1 a_2 \dots a_{K-1}$ . They are basically the elements of  $t_i$ , this is the Eigen vector associated with  $\lambda_i$  and  $\lambda_i$  is very closed to 0, so  $t_i$  is the Eigen vector associated with  $\lambda_i$  and  $t_i$  is equal to  $a_1 a_2 \dots a_{K-1}$ .

So, maybe I will just give the little motivation behind this, you know if  $\lambda_i$  is closed to 0 then the condition index associated with  $\lambda_i$  is large, I mean that would be greater than 1000. And I mean and then you will get one linear dependents between the regressor variables associated with  $\lambda_i$ , so that is why you know as I mentioned

before also you know that the number of condition indices greater than 1000 that measures the number of linear dependencies in the data. And corresponding to each lambda i for which the condition index is greater than 1000, you will get a linear dependents or you will get you can identify linear dependents between the regressor coefficient.

(Refer Slide Time: 19:16)

Webster Data

The smallest eigen value is  $\lambda_6 = 0.0011$  ←

$$t_6 = \begin{pmatrix} -0.447 \\ -0.421 \\ -0.541 \\ -0.573 \\ -0.006 \\ -0.002 \end{pmatrix}$$

$$-0.447x_1 - 0.421x_2 - 0.541x_3 - 0.573x_4 - 0.006x_5 - 0.002x_6 = 0$$

$$\Rightarrow x_1 = -0.941x_2 - 1.21x_3 - 1.28x_4$$

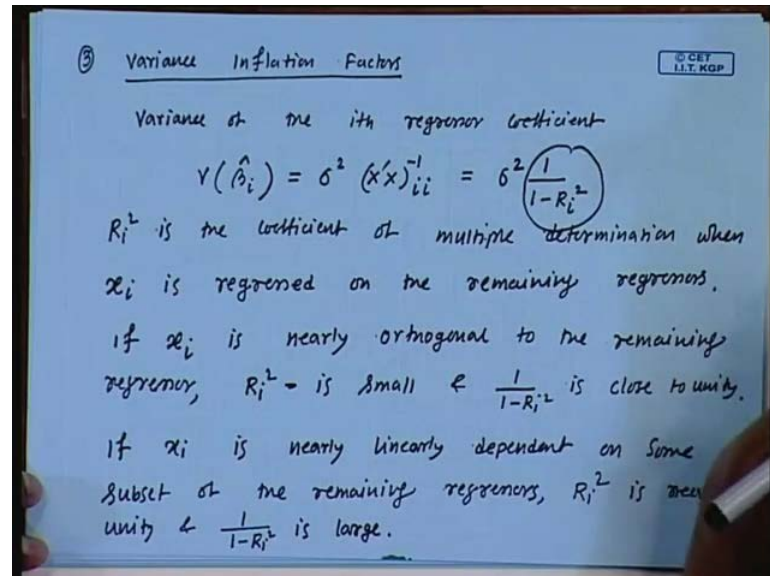
Let me explain I mean illustrate this thing using the webster data, here the smallest Eigen value is lambda 6 which is equal to 0.0011 and the associated Eigen vector that is say t 6 is equal to minus 0.447 minus 0.421 minus 0.541 minus 0.573 minus 0.006 minus 0.002. So, this is the associated Eigen vector corresponds to lambda 6 and then the nature of linear dependents is minus 0.447 x 1 minus 0.421 x 2 minus 0.541 x 3 minus 0.573 x 4 minus 0.006 x 5 minus 0.002 x 6 equal to 0.

And from here you know since these two are very small, we can ignore x 5 and x 6 from this equation and this implies that x 1 is equal to minus 0.941 x 2 minus 1.21 x 3 minus 1.28 x 4. So, this is the linear dependents between the regressor x 1 x 2 and x 3 and x 4, and this linear dependents is associated with lambda 6. And if you have more Eigen value which are closed to 0, corresponds to each lambda i mean which are very small or which are closed to 0, you will get a linear dependents like this. So, what we learned from this Eigen system analysis is that it can detect the presence of multicollinearity, it



can measure the number of linear dependencies in the data, and also it can identify the nature of linear dependencies in the data.

(Refer Slide Time: 23:16)



So, next we move to the variance inflation factor, this is another way to you know to detect the presence of multicollinearity, this is called variance inflation factors. So, first we will recall the variance of  $i$ th regressors coefficient, I mean variance of  $i$ th regressor coefficient means the variance of the least co estimate of  $i$ th regression coefficient  $\beta_i$ .

We know that this one is equal to  $\sigma^2 (X'X)^{-1}_{ii}$ , now you know this can be proved that this  $i$ th element is equal to  $\frac{1}{1-R_i^2}$   $\sigma^2$  here. Now, what is this  $R_i^2$   $R_i^2$  is the coefficient of multiple determination, when  $x_i$  is regressed on the remaining regressors.

Now, see if the  $i$ th regressor  $x_i$  is nearly orthogonal to the remaining regressor, here you know this nearly orthogonal I mean that if the  $i$ th regressor is independent of the remaining regressors. Then  $R_i^2$  is small and  $1 - R_i^2$  is close to unity well, so the meaning of this one is that no  $x_i$  is nearly orthogonal to the remaining regressors; that means,  $x_i$  is independent of the remaining regressors. That means, there is no linear dependents associated with  $x_i$ , I mean  $x_i$  cannot be represent in terms of the linear combination of the other regressors.

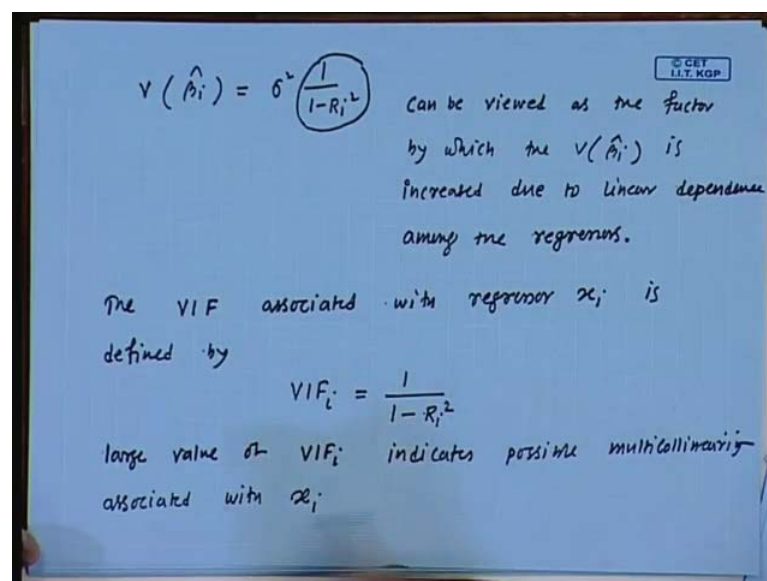


Then the coefficient of multiple determination when you are regressing  $x_i$  on the remaining regressors the coefficient multiple determination will be small, and the value of  $1 - R_i^2$  is close to unity, and the variance of  $\hat{\beta}_i$  is going to be  $\sigma^2$  close to  $\sigma^2$ . Now, if  $x_i$  is nearly linearly dependent on some subset of the remaining regressors,  $R_i^2$  value will be near to 1,  $R_i^2$  is near unity and  $1 - R_i^2$  is large.

So, the meaning of this, you know  $x_i$  is linearly dependent on some subset of the remaining regressors; that means, there is a linear dependence between  $x_i$  and some subset of the remaining regressors. If some linear dependence is there between  $x_i$  and some subset of the remaining regressors; that means,  $x_i$  can be represented in terms of as a linear combination of some subset of the remaining regressors well.

So, which implies that  $R_i^2$  which is the coefficient of multiple determination when  $x_i$  is regressed on the remaining regressors is will be large and that will be close to unity. And that implies that the value of  $1 - R_i^2$  is large, that ultimately you know if  $x_i$  is there is a linear dependence in the data between the regressors. Then variance of  $\hat{\beta}_i$  is going to be large because this factor is going to be large well, so this factor you know this  $1 - R_i^2$  well let me write here.

(Refer Slide Time: 30:48)

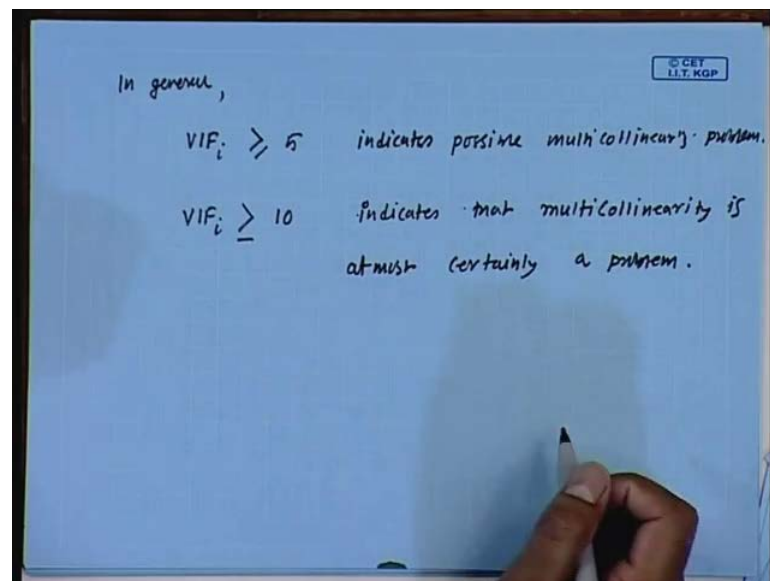


So, variance of  $\hat{\beta}_i$  is equal to  $\sigma^2$  into  $1 - R_i^2$  and this quantity this factor you know this can be viewed as the factor, by which the variance of

$\hat{\beta}_i$  is increased due to linear dependence among the regressors well. So, if there is a linear dependence between  $x_i$  and a subset of the remaining regressors, then the value of this is large, then the variance of  $\hat{\beta}_i$  is also large. Now, if  $x_i$  is independent of the remaining regressors or I say that it is if it is nearly orthogonal to the remaining regressors, then this value this factor is closed to 1 and the variance of  $\hat{\beta}_i$  is almost equal to  $\sigma^2$ .

So, the variance inflation factor VIF associated with regressor  $x_i$  is defined by  $VIF_i$  which is equal to  $1 / (1 - R_i^2)$  and; obviously, large value of  $VIF_i$  indicates possible multicollinearity associated with  $x_i$ . The meaning of this one is that you know if multicollinearity associated with  $x_i$ ; that means, this we if this is large; that means, there is a linear dependence between  $x_i$  and subset of the remaining regressors. Then only the value of this one is going to be large, so the large value of  $VIF_i$  indicates possible multicollinearity associated with the regressor  $x_i$ .

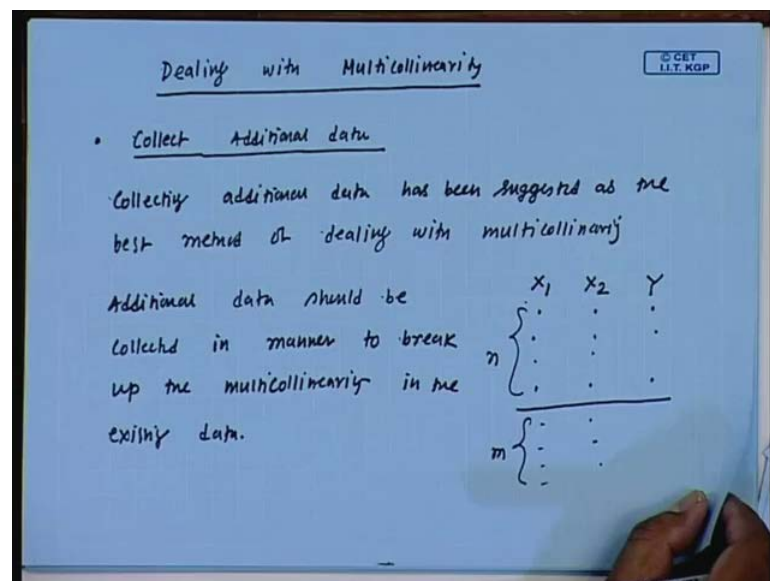
(Refer Slide Time: 34:50)



Now in general you know  $VIF_i$  greater than equal to 5 indicates possible multicollinearity problem, and  $VIF_i$  greater than or equal to 10 indicates that multicollinearity is almost certainly a problem. So, this is how the variance inflation factor associated with the  $i$ th regressor can be used to detect the multicollinearity, so this technique variance inflation factor can determine can detect the problem of multicollinearity.

But, any I mean of course, it cannot identify the nature of multicollinearity, so Eigen system analysis is a better technique, because it can detect the multicollinearity, it can measure the number of linear dependencies in the data, and also it can identify the nature of linear dependences in the data. So, next we will be talking about you know if you can if you detect that the you know there is multicollinearity in the data, then how to deal with multicollinearity, so we will be talking about several techniques to deal with multicollinearity.

(Refer Slide Time: 37:28)



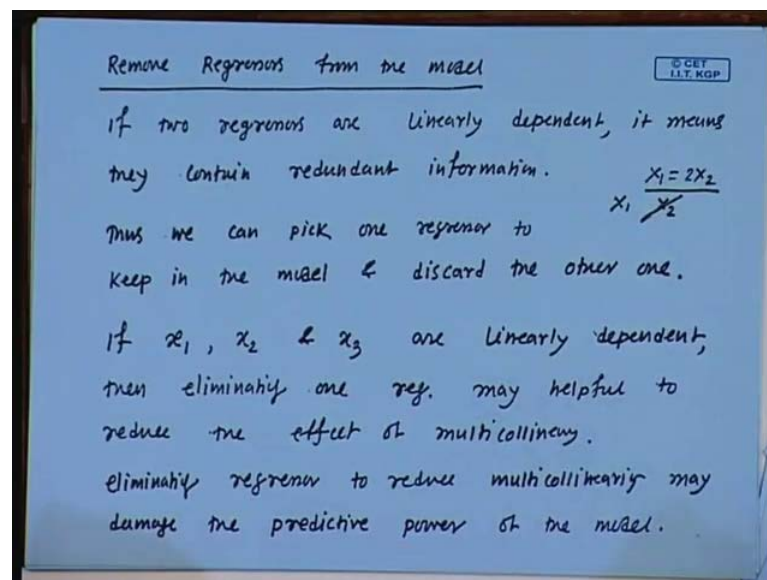
So, dealing with, so first technique is you know collect additional data, well you know collecting traditional data has been suggested, as the best method of dealing with multicollinearity. What it says is that let me illustrate this one it will get ,suppose you have you are in the multiple linear regression model, and suppose you have only two regressors  $x_1$  and  $x_2$  and the response variable  $Y$ , and you have end data points.

So, you have some data you have end data, now you have detected that you know multicollinearity exist in this data; that means, here we have only two regressors; that means,  $x_1$  and  $x_2$  they are linearly dependent. So, what we do is that we collect some more data say another  $m$  data to break the existing multicollinearity in the present data. So, what we do is that suppose, so this two regressors they are linearly dependent that is why that is the I mean multicollinearity, because of the linear dependence between these two regressors.

Now, this additional data should be collected in a manner to break up the multicollinearity in the existing data, I hope you understood that you know initially you had  $n$  data points, and here you know  $x_1$  and  $x_2$  they are linearly dependent. So, you collect another some more additional data say  $m$  data points, in such a way you know, when you combine the complete set of data you know  $n + m$  data then  $x$  and  $x_1$  and  $x_2$  are not any more linearly dependent.

You have to collect the data in a manner to break up the multicollinearity in the existing data, so this is you know one way to deal with multicollinearity or to break the multicollinearity in the existing data, but in many instance, you know this is not possible in practice.

(Refer Slide Time: 42:38)



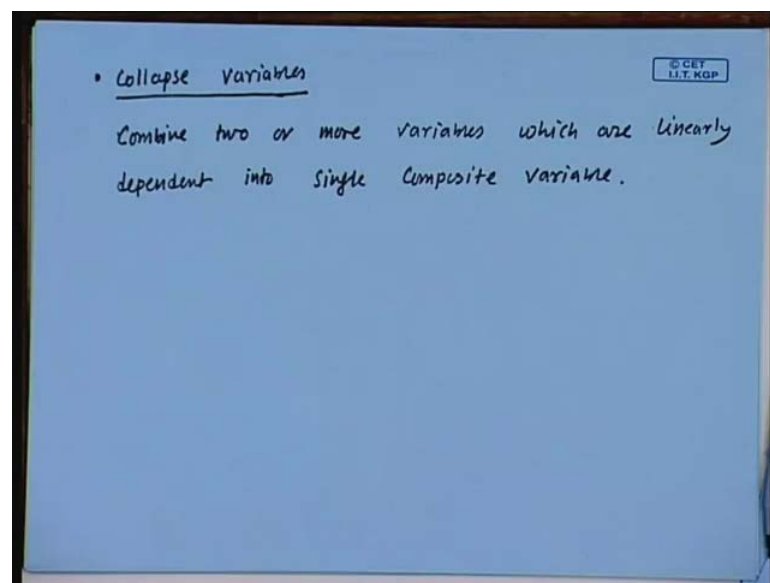
Next, we will talk about one more technique to deal with multicollinearity that is called remove regressors from the model, here if two regressors are linearly dependent, it means the content redundant information. So, what we can do is that you know if two regressors for example,  $x_1$  and  $x_2$  they are linearly dependent may be for example, say  $x_1$  is equal to twice of  $x_2$ , but then I mean basically they I mean the information regarding  $x_2$  is redundant.

So, what we can do is that thus we can pick one regressor to keep in the model and discard the other one, so this basically suggest you know, if you have two regressors, which are linearly dependent, and you can remove the other one. You can remove say for

example,  $x_2$  from this model, if say for example,  $x_1$ ,  $x_2$  and  $x_3$  are linearly dependent, then eliminating one regressor variable, may helpful to reduce the effect of multicollinearity, but the problem is that you know if say for example, you have three regressors,  $x_1$ ,  $x_2$ ,  $x_3$  in the model.

And or may be more regressors, but  $x_1$ ,  $x_2$ ,  $x_3$ , they are linearly independent, then maybe you can remove one regressor for example,  $x_3$  from the model, but you know to reduce the effect of multicollinearity. But, if  $x_3$  might happen that you know  $x_3$  the regressor which you have removed that  $x_3$  might be significant to explain the variability in the response variable. In that case you know this removing one regressor may damage the predictive power of the model, well so that is why it says that you know eliminating regressor to reduce multicollinearity may damage the predictive power of the model.

(Refer Slide Time: 48:20)



So, this is one way you know to deal with the problem of multicollinearity, the other technique is collapse variables, so it says that you know you combine if there are linear dependence between two or more than two regressors, which are linearly dependent into a single composite variable. So, basically this collapse variable it says that if you have linear dependence between say 2 or more than 2 regressors, then you can combine you know those regressors by a composite regressor variable.

So, these are the techniques to deal with multicollinearity, so one is you know collecting more data point, and the other one is removing one regressor from the model and

combining the regressors, which are linearly dependent. So, that is like ways for in this module we have learned what is multicollinearity and you know this multicollinearity is the name of a problem in multiple regression model well.

So, the problem of multicollinearity arises, when two or more regressor variables are linearly dependent and we have learned you know how to detect the multicollinearity, if it exist in the data. And also first you know we learned about how what are the problems due to multicollinearity, and then we have learned you know how to detect multicollinearity, if it exist in the data and also we learned how to deal with multicollinearity, so that is all for today.

Thank you for your attention.