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

**Lecture - 15**
**Multicollinearity (Contd.)**

Hi, this is my second lecture on Multicollinearity. And in the previous lecture, we have learned, what is multicollinearity? The problem of multicollinearity arises, when 2 or more than 2 regressors variables or linearly dependant. Today basically you know, we are in a last class also, we talked about the presence of multicollinearity has severe effect on least square estimates of regression coefficients. Today again, you know we will talk about the effects of multicollinearity or the problems due to multicollinearity and also we will learn how to detect the presence of multicollinearity and the data.

(Refer Slide Time: 02:08)



So, let me first talk about know, some more problems, due to multicollinearity, last class we proved or we illustrated the fact that, you know the strong multicollinearity results in large variants and co variants of the regression coefficients. So, we illustrated this issue using you know, in the case of multiple linear regression use, when there are 2 regressors in the multiple linear regression model, in also in general, we have proved this fact that, the strong multicollinearity results, in large variants and co variants of regression coefficients.

And today, we will be talking about some more problems, due to multicollinearity, the second one is multicollinearity tends to produce least square estimate beta hat that, are too far from the true parameter beta. So here, we compute the square distance between the least square estimate, of the regression coefficient, beta i hat minus beta i. So, this is the estimate of the i th regression coefficient and this is the true value of the i th regression coefficient, in the square distance is this one and sum over i is from 1 2 K minus 1 right.

And we denote, this by L square, next what we do is the, this is the square distance from beta hat to the true parameter value beta. Next, we compute the expected value of this square distance, so expected square distance, we denote it no this is basically, expectation of L square, which is equal to expectation of beta i hat minus beta i square sum over all i right well. What you know is that know, this beta i hat is the least square estimate of beta i and we know that the least square estimates unbiased estimated here, so we know that, the expected value of beta i hat is equal to beta i.

So, this one is nothing but expected value of beta i hat minus expectation of beta i hat right and this one is nothing but the variants of beta i hat. Now, the variants of beta i hat in multiple linear regression model is equal to summation sigma square X prime X inverse the i i th element. And this one is equal to sigma square by the i ith element, which we proved before also, it is 1 minus R i square sum over i, this one is basically sigma square sum 1 by 1 minus R i square. So, this R i square is the coefficient of multiple determination square, where R i square is the coefficient of multiple determination, for the regression of x i on the remaining K minus 2 regressors.

(Refer Slide Time: 09:51)



And also we, so what we proved is that, expected value of this square distance E L square is equal to sigma square sum 1 by 1 minus R i square sum over i and when, there is multicollinearity 1 minus 1 1 by 1 minus R i square will be large, for at least 1 i well. So, this R i square is the is the coefficient of multiple determination, for the regression of x i on the remaining K minus 2 regressors, now you know multicollinearity means, the problem of multicollinearity arises, when 2 or more regressors are linearly dependant.

Now, if say for example, if the i th regressors x i is linearly dependant on the remaining regressors then R i square that coefficients of multiple determination associated with x i is close to unity that means, you know R i square will be will tend to 1. When R i square will tend to 1 then 1 by 1 minus R i square will tend to infinity and then the expected value will be, I mean that is why it says, when there is multicollinearity, this term will be large for at least 1 i. So, it depends on you know, if x i is linearly depended on the remaining regressors than, this term will be large that and this term means large means, R i is close to I mean, when well, 1 minus 1 sorry, 1 by 1 minus R i square will tend to infinity as R i square tends to 1 well.

(Refer Slide Time: 13:19)



So, next we will talk about some more problems, due to multicollinearity that is the model coefficient with negative sign, when positive sign is expected well. So, that it says that, you know, you may get negative sign, for some regression coefficient, when you expect, you know, you really expecting positive sign for that regression coefficients. So, this might be the effect of multicollinearity and the next issue is say 4 says that high significance in a global F test, but in which none of the regressor are significant in partial F test.

So, I like to illustrate this point also, it says that you know high significance in the global F test, but none of the regressors are significant in partial F test. So, for this one, I want to recall one example, from multiple linear model, I mean there is only 1 example in I mean, I talked about in model 2 that means, in multiple linear regression model, I am going to recall that example. So, here is the data, I mean or here is the example.

(Refer Slide Time: 16:13)



We considered in model 2 multiple linear regression model well, so here, we have 2 regressors and 1 responds variable and we have the data for that well. We know how to fit a multiple linear regression model, here the fitted model, if you can recall the fitted model is Y hat equal to 14 minus 2 X 1 minus X 2 by 2. Now, once we have the fitted model, what we do is that, we check the significance of the fitted model by using the global test, that means, we test the hypothesis that, beta 1 equal to beta 2 equal to 0, against the alternative hypothesis that, you know H naught is not true. So, the significance of the null hypothesis that, beta 1 equal to beta 2 equal to 0, it says that you know, there is no linear relationship between the responds variable and the regressor variables. And we test this hypothesis, using the global F value that is obtained from the ANOVA table.

(Refer Slide Time: 18:11)



So, here is my ANOVA table, for this data, just refer my classes in model 2 well, so we have perhaps 11 data, this that is why the total degree of freedom is 10 and the F value is 7.17 and this follows, this F follows F distribution with degree of freedom 2 and 8 right. And since the observed F value is larger than, I mean greater than the tabulated F value, we reject the null hypothesis, beta 1 equal to beta 2 equal to 0 and that means, we accept, the alternative hypothesis that is beta 1 beta i is not equal to 0, for at least one i.

So, the ultimate conclusion from this test is that you know, the fitted model is significance the global test, says you reject the null hypothesis that, the null hypothesis is you know, it is it says that, there is no linear relationship between the responds variable and the regressor variable. So, we are rejecting that null hypothesis, that means, we are accepting the fact that, there is linear relationship between the responds variable and the regressors variables, now we go for the partial F test.

(Refer Slide Time: 20:16)



So, here is my partial F test, the first it says that, you know, what does X 2 contributes given that, X 1 is already in the regression. So, the contribution of X 2 in the presence of X 1 is you know, whether the contribution is significant that, can be tested by testing, the hypothesis that, beta 2 equal to 0 against beta 2 is not equal to 0. So, this 1, basically it test the significance of X 2 in the presence of X 1 in the regression model and you know either, you can go for the partial F test or you can go for the t test to test, this hypothesis well.

So, here I took the t statistic approach, here the t value is this one, now the tabulated t value is 2.306 right and the observed value is not greater than the tabulated value, so that means, we accept the null hypothesis that, beta 2 is equal to 0. So, accepting this null hypothesis means, accepting the null hypothesis beta 2 equal to 0 means that, the regresssor, second regressor X 2 is not significant in the presence of X 1. Now, what we do next, what will do is that, we will test the significance of X 1, now in the presence of X 2 well.

The next one is what does X 1 contributes given that, X 1 is already in the regression that means, the significance of X 1 in the presence of X 2 in the model, this can be tested by testing the hypothesis that, you know H naught equal H naught beta 1 equal to 0 against beta 1 naught equal to 0. Here is the test statistic value, here is the tabulated value, again you see that, the tabulated value sorry, the observed value is not greater than the

tabulated value that is why, we accept the null hypothesis beta 1 equal to 0. So, that means, none of the partial F test are significance partial test or the t test none of the partial test are significance.

(Refer Slide Time: 23:31)



So, what we observed here is that see, the global F test is significant the global F test says that, there is linear relationship between the responds variable and the regressor variable.

(Refer Slide Time: 23:47)

But when, we go for the partial test than, none of the regressors are significant neither X 2 is significant in the presence of X 1 nor X 1 is significance in the presence of X 2. So, this is one example basically, if the global F test is significant than, at least there should be one regressor, which is significant. But, here we are getting the result know the global F test is significant, but none of the partial test are significant; that means, this is this might be the effect of multicollinearity. So, you can check, you know whether multicollinearity exists in the given data that means, whether X 1 and X 2 are linearly depended that, you can check. So, this is one example of the effect of multicollinearity in the data right.
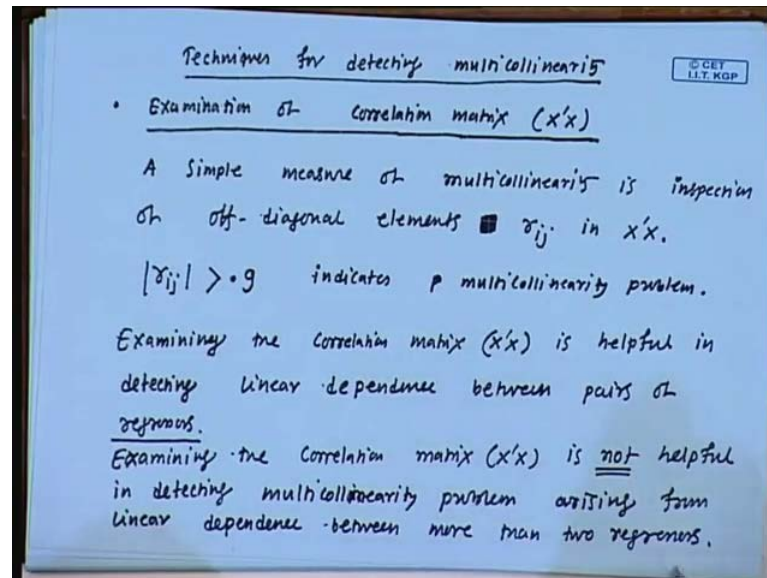
(Refer Slide Time: 24:55)



So, next we moved to another effect of multicollinearity that is that says that, you know different model selection procedures yield different models ok. So, model selection means you know about, the model selection, we talked about selecting the best model, in model 3 perhaps well. We know, how to select the best model, using the all possible selection and also the stepwise selection well, what it says that, if you have, if there is the presence of multicollinearity in the data then different model selection techniques will produce different model.

So, this is if you know a course of a different model can, different model selection procedure can yield different models, but if there is multicollinearity present in the data then the different with high probability that know, the different model selection

procedure will yield different models. These are the you know the different problem that, can occur due to multicollinearity, we talked about, next we will be talking about the I mean, we will talk about, the different techniques to detect the multicollinearity ok.

(Refer Slide Time: 27:26)



The first technique is examination of correlation matrix X prime X, so it says that, simple measure of multicollinearity is of off diagonal elements that is r i j in X prime X. So, we know, what is correlation matrix given the original data, if you scale and center them then the transformed data or the modified data, the X prime X for the modified data is called the correlation matrix.

Now, the off diagonal elements of the correlation matrix is are r i j, where r i j is the sample is the correlation coefficient between, the regressor x i and x j, now you know, if regressor x i and x j, if there dependant then the r i j will be near to unity. So, that means, if the regressor x i and x j, the correlation coefficient for the x i and x j value is high that, indicates the presence of multicollinearity. So, as a general rule, we say that you know, if r i j r i j is the correlation coefficient between the regressor x i and x j, if this one is greater than 0.9 then it indicates r i j greater than 0.9 indicates multicollinearity problem.

Now, examining the correlation matrix X prime X is helpful in detecting linear dependence between pairs of regressors, but the same is not no let me complete you now. But, examining the correlation matrix X prime X is not helpful in detecting multicollinearity problem arising, from linear dependence between more than 2

regressors. So, let me explain, what I wrote here, it says that the examining the correlation matrix is helpful in detecting linear dependence between pairs of regressors.

But, the same is not helpful in detecting multicollinearity problem arising from linear dependence between more than 2 regressors well. So, what it say basically, you know, if the problem of multicollinearity is due to the linear dependence between 2 regressors than, the correlation matrix can detect that, if the multicollinearity problem is due to the linear dependence between 2 regressors.

But, it might be the cases, we are talk, we are in the multiple linear regression setups, so there are K minus 1 regressors and K minus 1 could be greater than 2, but it might be the case that, you know the problem of multicollinearity is due to the linear dependence between more than 2 regressors. In that case, this correlation matrix cannot detect that, presence of multicollinearity in that case well. So, let me give one example to illustrate this fact, I will be taking this data.
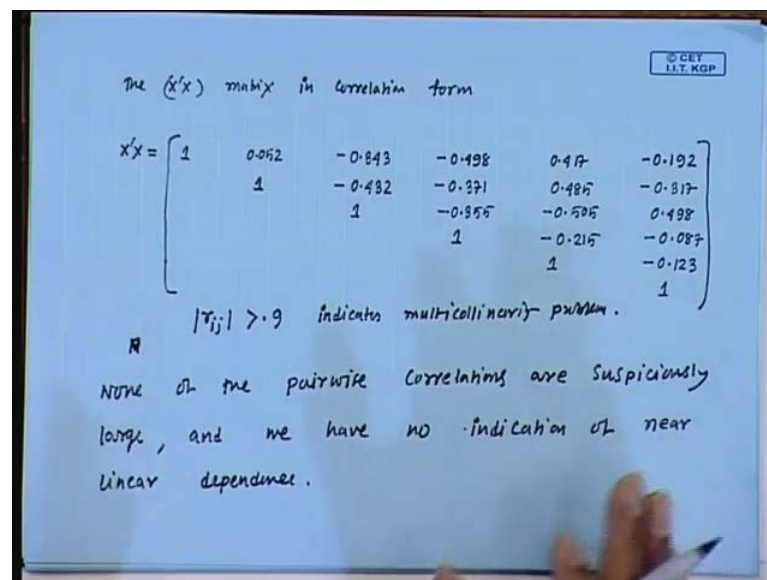
(Refer Slide Time: 36:35)



It says that, you know unstandardized regressor response variables, from webster gunst and mason well. So, I will refer this data, as webster data, here we have 6 regressors response variable and to check you know, what I want to say here is that, latter on we will see that, you know this data, I mean this data has the problem of multicollinearity, but the correlation matrix cannot detect it.

Because, the multicollinearity problem in this webster data is not due to the linear dependence between 2 regressors, here we have linear dependence between the regressors involving more than 2 regressors. So, what we do is that, first we will compute, you know first, we will compute the correlation matrix, for this data.
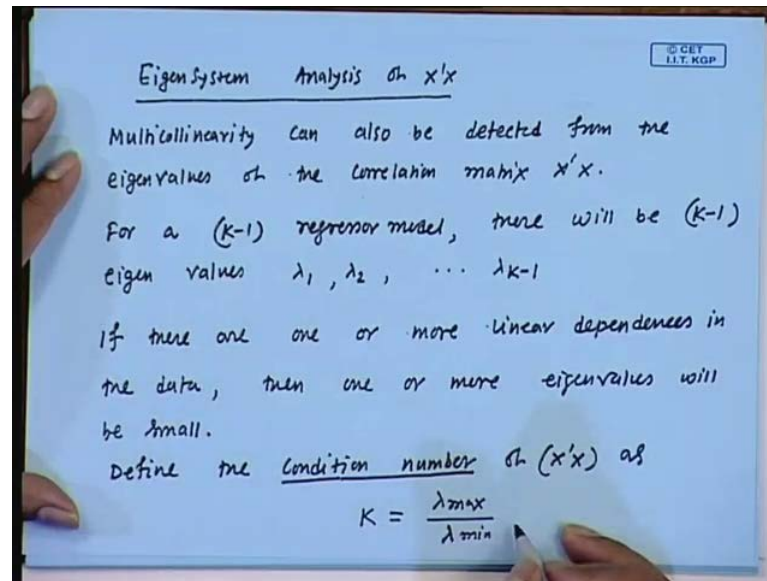
(Refer Slide Time: 38:07)



And here is the correlation matrix for the webster data, you see the off diagonal elements here, none of them are suspiciously large, here you know this is perhaps the highest correlation value. So here, we say, you know the we say that, r i j, so this is r 1 2, this is r 1 3, we say that, you know r i j, if r i j is greater than 0.9 that indicates, this indicates multicollinearity problem. But, here none of the pair wise correlations are suspiciously large and we have no indication of near linear dependence.

So, here you know the inspection of r i j is not sufficient to detect the multicollinearity, so the ultimate conclusion is that you know, the examining the correlation matrix is not sufficient, to detect the multicollinearity problem, if the multicollinearity problem involves linear dependence between more than 2 regressors well. So, next we will be talking about, one more techniques to detect the multicollinearity, so the examining the correlation matrix is not sufficient always.
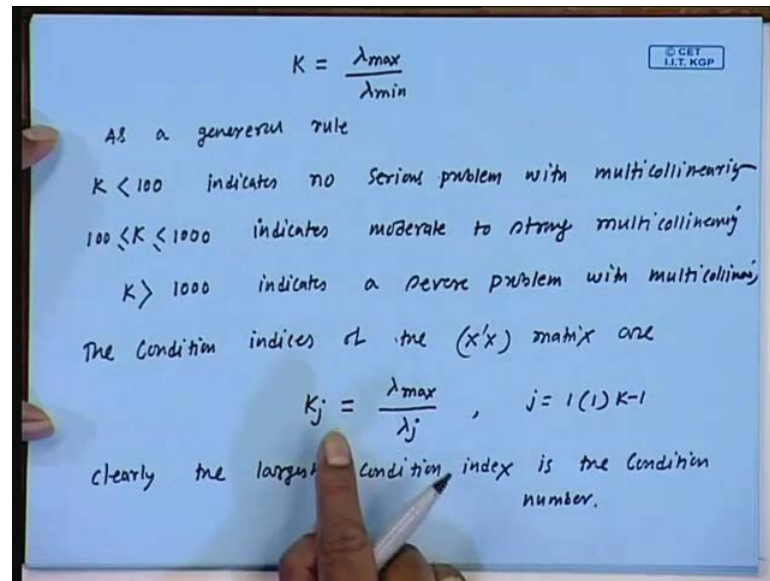
(Refer Slide Time: 41:00)



So, next we talk about, Eigen system analysis of X prime X, so it says that, multicollinearity can also be detected from the Eigen values of the correlation matrix X prine X. So, I hope you know, what is Eigen value and Eigen vector associate Eigen vectors, so for a K minus 1 regressor model, this one is K minus 1, cross K minus 1 matrix, there will be K minus 1, Eigen values say lambda 1 lambda 2 lambda K minus 1.

So, now, you know, if there are one or more linear dependence in the data then one or more Eigen values will be small well. So, basically you know, if you have a small Eigen values that implies, small Eigen values implies that, there is linear dependence between the columns of X. Now, we define the condition number of X prime X as K, which is equal to lambda max by lambda minimum, lambda max is the maximum Eigen value and lambda minimum is the minimum Eigen value.

Now, you know that, a small Eigen value indicates linear, I mean 1 or more small Eigen values indicates 1 or more linear dependences in among the regressor variables or linear dependences among the columns of X. Now, look at the condition number here, if lambda minimum is small or very close to 0 than this condition number is going to be large and since lambda minimum, if lambda minimum is small then you know there is linear dependences in the data. So, from there, we can conclude that, when lambda minimum is small, K is going to be large. So, the large value of K indicates the presence of multicollinearity in the data right.
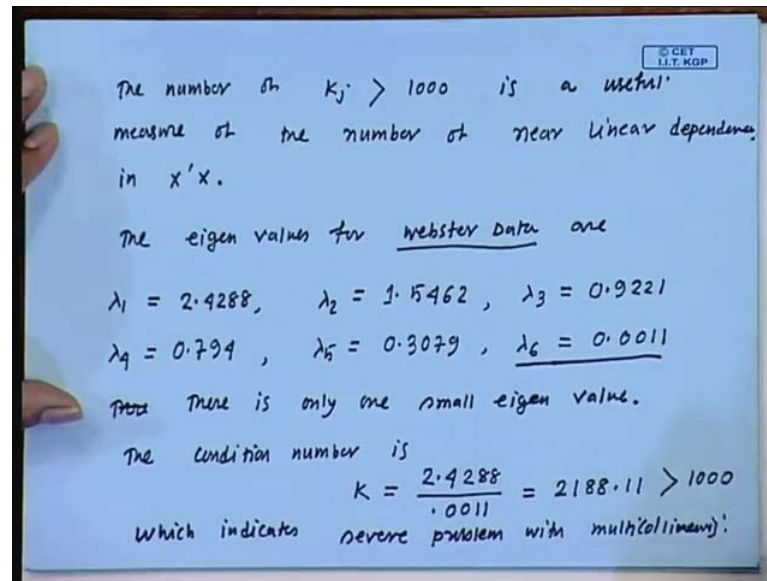
So, we will give a general rule, what is K, K is lambda max by lambda min, as a general rule K less than 100 indicates no serious problem with multicollinearity, now K in between 100 to 1000 indicates moderate to strong multicollinearity and K greater than 1000 indicates a severe problem with multicollinearity. So, large value of K indicates the severe problem with multicollinearity, because of the fact that, you know K will be large, when lambda minimum is very small or close to 0 well.

And since lambda minimum is 0 indicates, there is linear dependence among the columns of X, that means, the presence of multicollinearity. So, this is the condition number, now we define the condition indices of the X prime X matrix or K j, which is equal to lambda max by lambda j, for j equal to 1 to K minus 1 and of course, you know clearly the largest condition index is the condition number.

So, this will be large, when lambda j is minimum and that is nothing but the condition number, now this Eigen system analysis, it not only you know, it not only detect the multicollinearity problem also, it can measure the number of linear dependences in the data. So, that can be, I mean the number of K j or condition indices greater than 1000 is a measure of the number of linear dependences in the data.

(Refer Slide Time: 51:14)



So, the number of K j greater than 1000, because if see, if K j is greater than 1000 than that indicates severe problem with multicollinearity, now number of K j greater 1000 is a useful measure of the number of near linear dependences in X prime X. If we considered the webster data, the Eigen values for webster data are, so there, we had you know 6 regressor. So, 6 Eigen values lambda 1 equal to 2.4288, lambda 2 equal to 1.5462, lambda 3 equal to 0.9221, lambda 4 equal to 0.794, lambda 5 equal to 0.3079, lambda 6 equal to 0.0011.

So, this Eigen value the smallest Eigen value is very close to 0, so this one is basically you know, the there is we say that, there is only one small Eigen value means, you know very close to 0. And the condition number here, also you know small Eigen values indicates linear dependence in the data 1 linear, so since, there is only 1 small Eigen value, we can say that, there is there might be only one linear dependence in the data.

But of course, we have to check with the condition indices also, let me first compute the condition number is K, which is lambda, max lambda max is lambda 1 here 2.4288 by 0.0011, which is equal to 2188.11 and this is larger than 1000 right. So, this indicates, which indicates severe multi severe problem with multicollinearity, so what I want to say here today regarding the detection of multicollinearity problem is that, you know first we considered this webster data and examining, the correlation matrix would not detect the problem of multicollinearity in the webster data.

Where as, you know this Eigen system analysis technique, it says that since the condition number is 2188, which is much larger than 1 1000. So, it says that, the there is a severe problem with multicollinearity in the webster data, so not only, this one in the next, we will be talking more about, you know we will compute the condition indices and from there, we will see how many linear dependences are there. So, that is all for today.

Thank you for your attention.