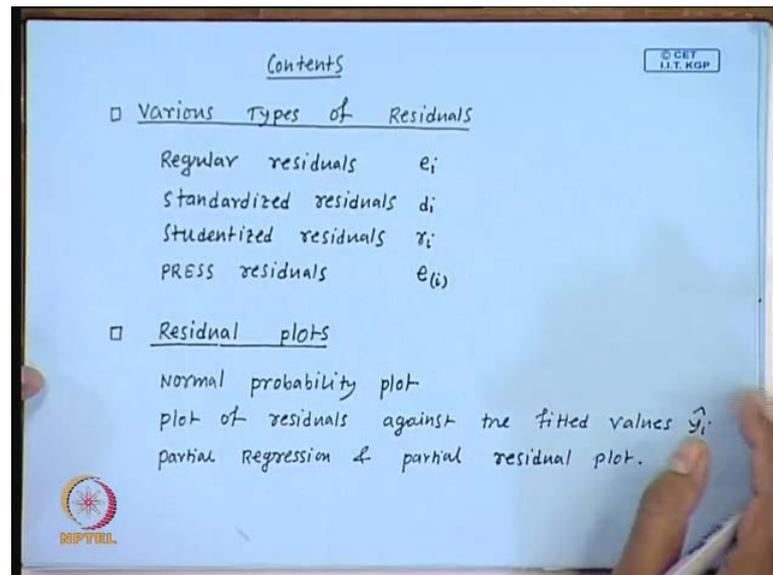


Regression Analysis
Prof. Soumen Maity
Department of Mathematics
Indian Institute of Technology, Kharagpur

Lecture - 18
Model Adequacy Checking (Contd.)

(Refer Slide Time: 00:28)

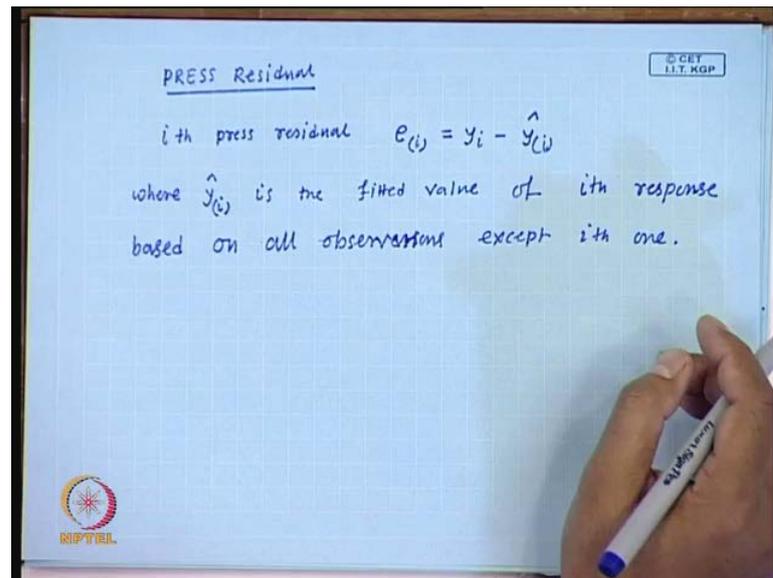


This is my second lecture in model 5, that is model adequacy checking and here is the content of this model various type of residuals. And in the previous class we talked about regular residuals, standardized residuals, studentized residuals. And today we will be taking about you know the press residuals, and also next we will be talking about several types of residual plots. Like, you know normal probability plot, plot of residuals against the fitted values \hat{y}_i and in the next class maybe we will be talking about partial regression and partial residual plot.

So, before I start talking about you know the press residual, just I want to repeat once more the objective of this model here, you know if you can recall you know simple linear regression model or in the multiple linear equation model. We have assumed that the error term ϵ has 0 mean, and error term ϵ has a constant variance, and the error terms are correlated and they are normally distributed. So, what we are doing going to do in this model is that, we will present you know several methods to check the underlying assumptions that we made on the error term ϵ .

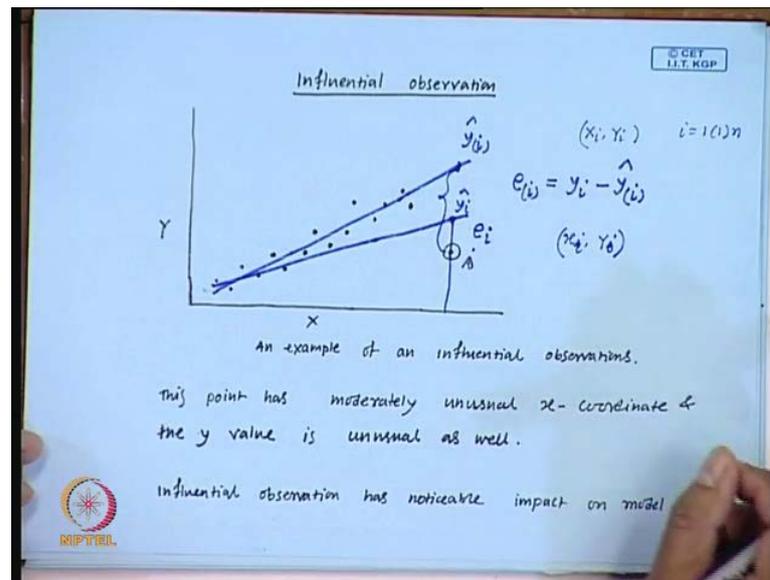
And the methods you know they mostly depend on the primarily depend on the study of residuals, because we think it is convenient to think that the residuals are the realized or observed value of the epsilon. Since, you are going to test some assumption on the error term epsilon, so the test for that is based on the residuals, and the graphical analysis of the residuals are very effective to test the underlying assumptions on epsilon.

(Refer Slide Time: 03:53)



So, now I will be talking about press residual, which is one scaled residual, so I have already talked about regular residual and the other things, today I will be talking about press residual. So, it is beautiful concept here, well the i th press residual denoted by $e_{(i)}$ is equal to y_i minus $\hat{y}_{(i)}$, so we know that y_i is the i th observed value, and what is $\hat{y}_{(i)}$ where $\hat{y}_{(i)}$ is the fitted value of i th response based on all observations except i th one. So, the basic logic behind this that, so this is not the regular residual this is the i th observation the response, value of the response variable and this is the fitted value of the i th response based on all observations except the i th observation.

(Refer Slide Time: 06:29)



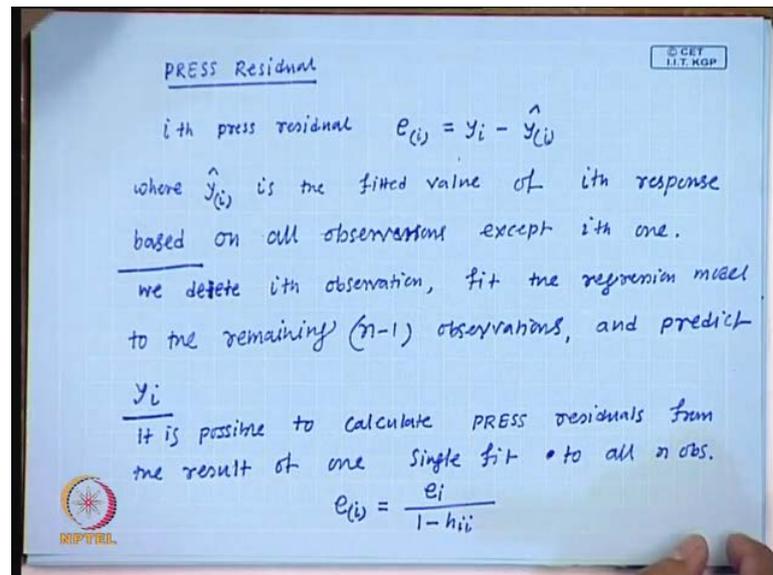
Let me just give idea about this special type of press residual, I will just refer the previous example, this is an example of influential observation. Now, if I fit a model you know based on the observation share, including this influential observation also y fitted model may get influenced by this influential observation. So, my fitted model may look like this, so this is the fitted model, because of this influential observation, and you can see you know the fitted model has been influenced by this influential observation.

And the fitted value here I mean this is the model based on all the observations for the fitted value here is y_i hat, so this is my y_i hat, and this is I am considering this as you know the i th observation. Then the residual I mean the true value of the i th observation, that is true value y_i and the y_i hat they are the difference is less. So, the difference is this much, and this is the i th regular residual e_i , now if this i th observation is deleted from this data set, suppose this observation is not there in the model in the set of observations.

Then my fitted model will look like this and this model this fitted model is of course, it is not influenced by this influential observation, because I am fitting this model based on all observations except this i th observation. So, this is what I am talking this value you know the this value is y_i hat bracket, so this is what I mean you know the fitted value of the i th response variable based on all the observation except the i th observation.

Now, my residual the reference residual which I call here e_i bracket is this difference, so e_i bracket is equal to y_i , so this is y_i and this is \hat{y}_i bracket here. So, y_i minus \hat{y}_i bracket, and you can see that the i th based residual is substantially larger than the regular e_i .

(Refer Slide Time: 10:42)



Well, so what we do is that, we delete we delete i th observation, fit the regression model to the remaining n minus 1 observations and predict y_i . So, here it may appear that you know to compute the first press residual that is e_1 in bracket. you need to fit model based on all observations except the first observation. So, that is how you get e_1 bracket that is a first press residual, again you know you do not know which one is the influential observation, so you have to repeat this process in time.

So, to get e_2 bracket I mean the second press residual again you have to fit a model to the based on all the observations except the second observations; that means, you know you need to fit model based on n minus 1 observations. That you have to repeat in times, so but what are we going to do, it can be prove that you do not need to repeat this process n times to get the n press residuals, it can be done you know based on one regression fit based on all the observations.

So, here is the technique; however, so it says that it is possible to calculate press residual from the result of one single fit to all n observations, it can be proved that you know e_i , i th press residual is equal to e_i by $1 - h_{ii}$. Just recall that you know this e_i is the

regular residual and h_{ii} is the i th diagonal element of the hat matrices, and this is how you know this quantity they are same, so this is how we calculate the press residuals e_1, e_2, \dots, e_n .

(Refer Slide Time: 15:25)

Observation No.	Delivery Time (y) (minutes)	No. of Cases x_1	Distance (Feet) x_2
1	16.68	7	560
2	11.50	3	220
3	12.03	3	390
4	14.88	4	80
5	13.25	6	150
6	18.11	7	330
7	8.00	2	110
8	17.83	7	210
→ 9	79.29	30	1960
10	21.50	5	605
11	40.33	16	688
12	21.40	10	215
13	13.50	4	255
14	19.75	6	462
15	29.00	9	498
16	29.00	10	776
17	15.35	6	200
18	19.00	7	132
19	9.50	3	36
20	35.10	17	770
21	17.90	10	140
22	52.32	26	810
23	18.75	8	450
24	19.83	8	635
25	10.75	4	150

© CEY
I.I.T. KGP

Ref:
Montgomery
Peck
Vining

$$\hat{y} = 2.34 + 1.615x_1 + 0.0143x_2$$

$$e_i = Y_i - \hat{Y}_i$$

Now, let me just recall the example, we considered in the last class, so here is an example of multiple linear regression with two regressors x_1 and x_2 , and here is this response variable Y . And we suspect that I mean it is likely that the ninth observation is an influential observation or at least it is an average point, because x_1 coordinate is much larger compare compared to the centre of x_1 and similarly x_2 value is much larger.

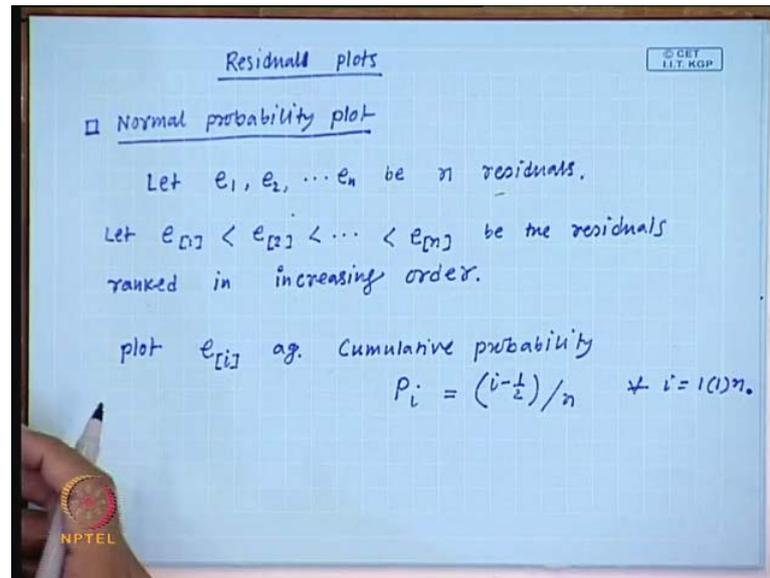
(Refer Slide Time: 16:09)

Observation (i)	$e_i = y_i - \hat{y}_i$	$d_i = \frac{e_i}{\sqrt{MS_{Res}}}$	$\hat{x}_i = \frac{e_i}{\sqrt{MS_{Res}(1-h_{ii})}}$	PRESS Residuals $e_{(i)}$
1	-5.0281	-1.54	-1.627	-5.895
2	1.1469	0.3517	0.349	1.233
3	-0.0498	-0.0143	-0.0161	-0.0557
4	4.9249	1.5108	1.579	5.2297
5	-0.444	-0.1363	-0.1418	-0.809
6	-0.2896	-0.0888	-0.0908	-0.3025
7	0.8446	0.2501	0.2709	0.9198
8	1.1546	0.3545	0.3667	1.2343
→ 9	7.4107	2.2743	2.2138	14.788 = $e_{(9)}$
10	2.3769	0.7291	0.8123	2.9568
11	2.2375	0.6865	0.7161	2.4989
12	-0.5930	-0.1819	-0.1932	-0.6690
13	1.027	0.3161	0.3252	1.0998
14	1.0675	0.3275	0.3411	1.1581
15	0.6712	0.2059	0.2103	0.700
16	-0.6629	-0.2039	-0.2277	-0.7948
17	0.4369	0.1339	0.1381	0.4640
18	3.4466	1.0580	1.1100	3.8159
19	1.7032	0.5402	0.5787	1.9846
20	-6.7880	-1.7758	-1.8736	-6.4932
21	-2.6142	-0.8020	-0.8779	-3.1318
→ 22	-3.6965	-1.131	-1.2500	-6.0491
23	-4.6076	-1.4136	-1.4437	-4.8089
24	-4.5723	-1.4029	-1.4961	-5.2000
25	-0.2126	-0.0652	-0.0674	-0.2278

And, here is the regular residual, this is standardized residual, this is the studentized residual, and this is the press residual. Now, here we check that for the ninth observation the value of the regular residual is 7.41, where as the value the ninth press residual, so this is nothing but this is equal to $e_{(9)}$ press residual, this is 14.788 whereas, the regular residual is 7.4. So, for the point for the ninth observation the press residual value is sustainably larger than the regular residual, now also we can observe that for the 20 second observation, the value of the regular residual is 3.6 whereas, the value of the press residual is minus 6.05.

So, here also the value of the press residual is substantially larger than the value of the regular residual for 22nd observation, so this is what regarding the press residual it appears that, if the i th observation is an influential observation or if it is leverage point. Then there will be substantial difference between the regular residual value and the press residual value for that particular observation. So, next we move to residual plots, so there are several residual plots, so we will be talking about the first one is normal probability plot.

(Refer Slide Time: 18:56)

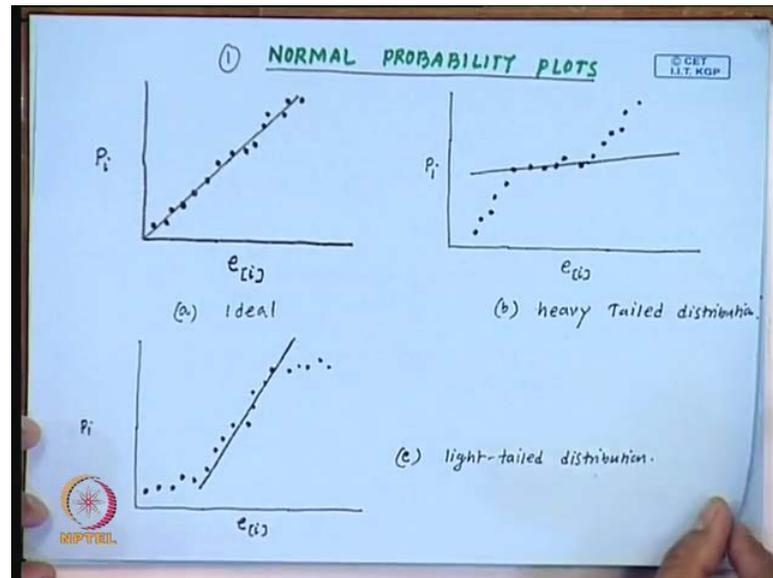


So, residual plots, first we will be talking about normal probability plot, so before I talk about this residual plot you know, I just want to mention again that the objective of this module is to check the underlining assumption on epsilon on the a rectum. So, we will be talking about different methods to check the underlining assumptions, so this is one of them. So, what this normal probability plot does is that it basically, check the assumption is that the epsilon i follows normal 0 sigma square, and they are independent and identical distributed.

So, this normal probability plot what it does is it checks whether the erectum really follow normal distribution or not, so here is the technique to test that. So, what it does is that let e_1, e_2, e_n be n residuals; that means, the regular residuals, let $e_{\text{box } 1}, e_{\text{box } 2}, e_{\text{box } n}$ be the residuals ranked in ranked in increasing order. So, given a set of data you know you can fit the model, whether it is simple linear regression or multiple linear regression.

And then you can get the residual values, you just rank them I mean you arrange them in increasing order. Then what this normal probability plot does is that it plots e_i against the cumulative probability p_i , which is equal to i minus half by n and you do it for all i , i is from 1 to n . So, very simple technique you get the residuals first you arrange them in increasing order, and then you plot e_i against p_i here is the p_i , i minus half by n .

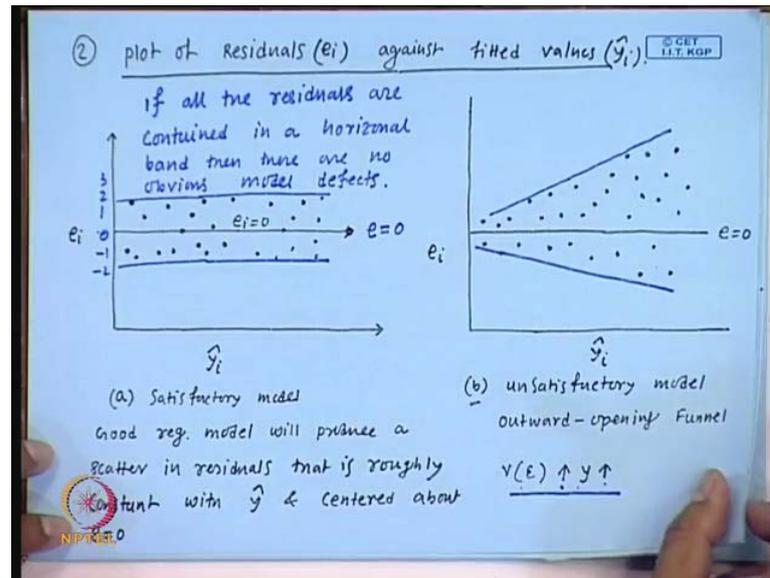
(Refer Slide Time: 23:36)



So, just now I have different types of normal probability plots, that may happen this is, so figure one shows the ideal situation, here you can see that all the points lie on straight line. So, if in your case, if you find that your plot I mean all the points are on the on a straight line, then you can assume that the error distribution is normal. Now, if here in the figure b, you can see that it is not really I mean the points are not on a single straight line, and here it this type of situation occurs.

Then we can say that the distribution is heavily tailed, that is you know it is really not true that the erect distribution is normal. And figure c indicates that the distribution of the erectum is lightly tailed. So, there is little deviation from the I mean it is not reasonable to assume, if such situation occur than it is not advice able to assume that the erect distribution, is normal. So, this is what regarding the residual plot, so residual plot basically, if test the normality assumption of the erectums epsilons. Next, we will be talking about one more plot, that is it is a plot of residual term against the fitted observation.

(Refer Slide Time: 26:26)



So, here is the plot of residual e_i against the fitted value \hat{y}_i , so given a set from data you know, you fit the model first and then you can get the fitted value. Once, you have the fitted value you can compute the residual, and then you plot the residual against the fitted value. Well now, we need to understand I mean of if the plot if the pattern looks like this what does it indicated, so if the plot of residual against the fitted value looks like this one.

Then you can conclude that the this is a good regression model, and the good regression model will produce scatter in residuals that is roughly constant with \hat{y} and centered about e equal to 0. Then I mean what I am trying to say here is that if all the residuals are contained in a horizontal band. So, here is the horizontal band, here you can see that all the points are contained inside this horizontal band, then there are no obvious model defects.

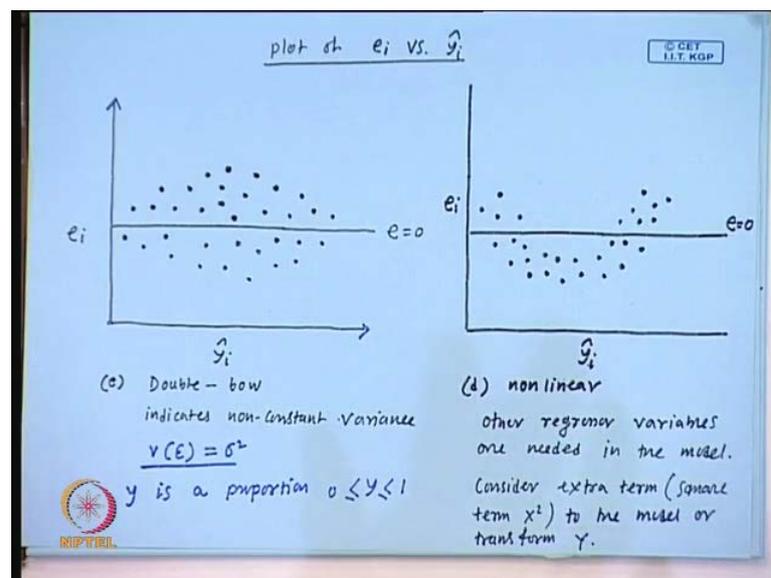
That means, if what I want to say here is that you know, if all the residuals are contained in a horizontal band centered about e equal to 0. This is the line e equal to 0, because residual value could be this is 0 minus 1 minus 2 1 2 3, like that some of the residuals are positive some are negative you have observed that before also. So, what I want to say is that if the residuals are contained in a horizontal band, then there is no obvious defect with the model; that means, it is a good fit.

So, this is the only situation we say that, the fitted model is satisfactory, so there is no problem with the fitted model. Now, look at the second case this is this figure b indicates, and how to outward opening funnel pattern, so here you can see that this e_i values increase, I mean e_i increases with the value with \hat{y}_i . So, this is called outward open funnel pattern, what does this indicates that the indicates non constant variants of it is, we assume the see we have assumed that the variants of the epsilon is constant that is sigma square.

But, if your residual I mean if your plot, which plot I mean this plot particular like epsilon verses epsilon sorry, e_i against \hat{y}_i . If the pattern is like you know outward open funnel, that indicates that variants of epsilon i or the variants of epsilon increases as a y increases. So, this one is sort of this is the indication of a non constant variants, so it cannot if this occurs you know it is not advisable to assume that the variants of epsilon is sigma square.

So, next we talk about some more pattern, the next pattern is one more thing just I forgot to mention here, instead of outward open funnel heat could be inward open funnel also. That means, it will look like this in that case also indicates that, non constant variants of epsilon, in that case variants of epsilon decreases as y increases. So, both the I mean it could be outward open funnel, it could be inward open funnel, so in both the case indicate the non constant variants of epsilon.

(Refer Slide Time: 33:30)

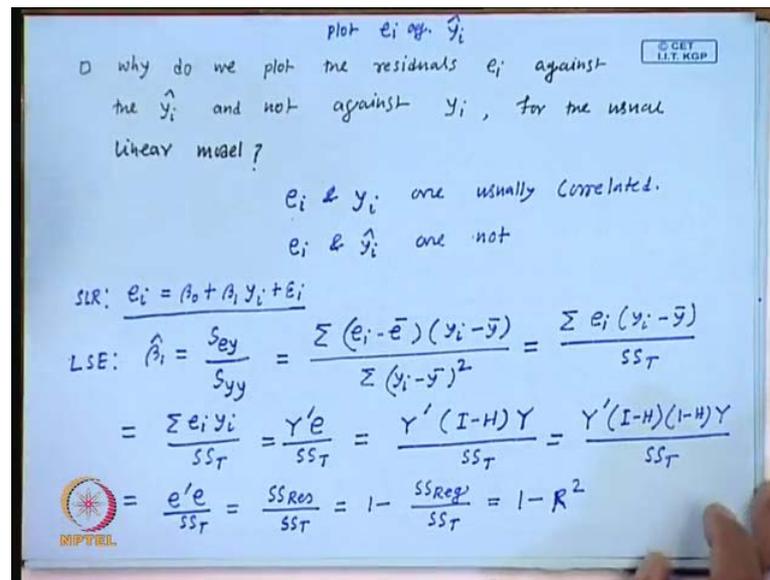


So, next we will be talking about one more I mean some more pattern, like a this is one more pattern here this is the line e equal to 0 y_i hat, in this direction, and this is called double bow. And this one also indicates non constant variants; that means, variants of epsilon, we cannot assume that this is equal to sigma square, so this pattern also violet this assumption. Well, and this type of pattern often occurs when y is a proportion, response variable is a proportion and y lies between 0 to 1, also this is in this type of situations we often get double bow pattern.

Now, the last one is that you know the figure d shows the non-linear pattern here, this non-linear pattern indicates that you know other regressor variables are needed in the model. So, this indicates that the relationship between y and the regressor variable is not linear, we need to introduce some non-linear term; that means, consider extra term like square term x square.

I mean the relation is not just the linear relation like y equal to beta naught plus beta 1 x you need to introduce some square term the higher order term, or you terms you take a transformation of the response variable y . I mean maybe we will talk regarding this issues later on again, but what this non-linear pattern indicates that the relationship between the response variable in the regressor variable is not linear. We need to introduce some other terms like, other regressors. For example, we need to take x square x^q in the model, or you may need to take a transformation of the response variable y like $\log y$ or something 1 by y something like that.

(Refer Slide Time: 36:43)



Now, there is a question you may be wondering you know why do, what we did is that we have just checked about another plot of e_i against \hat{y}_i , so why it is you know why plotting e_i against \hat{y}_i , why not e_i against y_i . So, this is the question why do we plot the residual e_i against \hat{y}_i and not against y_i for the usual linear model, it is not, so easy to answer this question like you can think of it, e_i vs y_i . The answer to this question is that you know e_i and y_i we do not constraint the plot e_i against y_i , because this two are usually correlated.

What I mean by this I will talk about that whereas, e_i and \hat{y}_i are not correlated, there is no relationship between e_i and \hat{y}_i , what I am going to prove is that there is a linear relationship between e_i and y_i . Suppose, the relationship is the form of the relation is e_i I said linear relationships e_i equal to some β_0 plus $\beta_1 y_i$ plus ϵ_i is just a simple linear relation between the residual and the observed value y_i .

Now, if you know the same technique this is nothing but the simple linear regression model between the residual and y_i , you can check that $\hat{\beta}_1$ is nothing but S_{ey} / S_{yy} . So, what is S_{ey} this is the list square estimate you can check with my first model $\hat{\beta}_1$ is nothing but this one this is called the list square estimate. So, this one is nothing but notation only this is nothing but $\sum (e_i - \bar{e})(y_i - \bar{y}) / \sum (y_i - \bar{y})^2$.

So, what am I trying to prove that there is a relationship between y_i and e_i of and I am trying to find out that relationship, what type of relationship they have, if it is linear then what is the value of the coefficient. So, this one is nothing but summation $e_i y_i$ minus y bar, you can check that it is not difficult and this one is nothing but the SST , and again you know y bar into e_i that is sum over e_i is going to be 0, so this is equal to summation $e_i y_i$ by SST very simple.

Now, in matrix notation this can be written as $e' y$ or $y' e$ same thing by SST , now what is e in terms of h rotation Y' and e we know e is $I - H$ Y , by SST . And we know what we know $I - H$ is idempotent matrix $I - H$ is also idempotent matrix, so $I - H$ can be replaced by $(I - H)^2$. So, this is equal to $Y' (I - H) Y$, because $(I - H) Y = Y - H Y = Y - e = e'$, as $I - H$ is idempotent matrix into Y by SST well, so this is nothing but see this is e' this is e prime this is e prime e by SST .

So, e' is e is nothing but SST residual this is nothing but SST residual by SST which is nothing but $1 - R^2$ regression by SST which is nothing but $1 - R^2$. So, this R^2 is the you can recall R^2 is the coefficient of multiple determination well, so the relationship between, so that what we proved is that there is a linear relationship between the residual and the observed value. And the coefficient value I mean the slope is equal to $1 - R^2$, well let me check with \hat{y}_i whether there is linear relationship between e_i and \hat{y}_i , we can prove that there is that the slope is 0 in that case.

(Refer Slide Time: 43:58)

$$e_i = y_i - \hat{y}_i$$

$$\text{SLR. } e_i = \beta_0 + \beta_1 \hat{y}_i + \epsilon$$

$$\text{LSF } \hat{\beta}_1 = \frac{S_{e\hat{y}}}{S_{\hat{y}\hat{y}}}$$

$$SS_{e\hat{y}} = \sum (e_i - \bar{e})(\hat{y}_i - \bar{y}) = \sum e_i \hat{y}_i$$

$$= e' \hat{y} = Y'(I-H)HY$$

$$= Y'(H-H^2)Y$$

$$= Y' \cdot 0 \cdot Y$$

$$= 0$$

$$e = (I-H)Y$$

$$\hat{y} = HY$$

$$H^2 = H$$

So, let me check with e_i and \hat{y}_i , so e_i suppose there is a linear relation between e_i and \hat{y}_i and the relation is e_i equal to $\beta_0 + \beta_1 \hat{y}_i + \epsilon$ like that the simple linear regression. And then by least square estimate we can check that $\hat{\beta}_1$ is equal to $S_{e\hat{y}} / S_{\hat{y}\hat{y}}$, I do not care about this denominator, let me proceed with $SS_{e\hat{y}}$.

If I can prove that this is equal to 0, then my slope is going to be equal to 0; that means, there is no linear relationship between e_i and \hat{y}_i . So, this is going to be equal to e_i minus \bar{e} into \hat{y}_i minus \bar{y} , and you can check that this is nothing but summation $e_i \hat{y}_i$. So, in matrix notation this is equal to $e' \hat{y}$, so what is e' , e' is equal to $Y'(I-H)$, because e equal to $(I-H)Y$, Y and \hat{y} is equal to HY . So, HY here now $Y'(I-H)HY$, now see H is idempotent matrix, so H^2 is going to be equal to H , so this is $Y'(H-H^2)Y$, this is going to be equal to 0. So, there is no, so this proves that there is no linear relationship if between the residual and the fitted values.

(Refer Slide Time: 46:54)

plot e_i vs. \hat{y}_i

□ why do we plot the residuals e_i against the \hat{y}_i and not against y_i , for the usual linear model?

e_i & y_i are usually correlated.
 e_i & \hat{y}_i are not

SLR: $e_i = \beta_0 + \beta_1 y_i + \epsilon_i$

LSE: $\hat{\beta}_1 = \frac{S_{ey}}{S_{yy}} = \frac{\sum (e_i - \bar{e})(y_i - \bar{y})}{\sum (y_i - \bar{y})^2} = \frac{\sum e_i (y_i - \bar{y})}{SS_T}$

$= \frac{\sum e_i y_i}{SS_T} = \frac{Y'e}{SS_T} = \frac{Y'(I-H)Y}{SS_T} = \frac{Y'(I-H)(I-H)Y}{SS_T}$

$= \frac{e'e}{SS_T} = \frac{SS_{Res}}{SS_T} = 1 - \frac{SS_{Reg}}{SS_T} = 1 - R^2$

Well so what I want to conclude here is that, unless in case of e_i and y_i we have observed that $\hat{\beta}_1$ is equal to $1 - R^2$.

(Refer Slide Time: 47:04)

e_i vs. \hat{y}_i , $\hat{\beta}_1 = 1 - R^2$

unless $R^2 = 1$, there will be a slope of $(1 - R^2)$ in e_i vs. \hat{y}_i plot, even if there is nothing wrong with the model.

So, in the case of e_i and y_i $\hat{\beta}_1$ is equal to $1 - R^2$, so unless R^2 is equal to 1, there is a positive slope, there will be a slope of $1 - R^2$, in e_i versus y_i plot. Even if there is nothing wrong, with the model, so if you can recall that in e_i

verses in e_i against \hat{y}_i plot, I said that if all the residuals are contained in a horizontal band, centered around e equal to 0.

Then the corresponding fitted model is perfect, but here since there is a theoretical relationship between the residual and \hat{y}_i it is very likely that the residuals will not be contained within horizontal band centered at E mean e equal to 0. There will be slope of one minus r^2 when r^2 is not equal to 0, so that is why you know it is very difficult to conclude anything, if we plot the residual against \hat{y}_i instead of plotting the residual against y_i . So, that is why the reason you know, and since there is no relationship between e_i and \hat{y}_i , they are not correlated there is no linear relationship between them, that is why we prefer plotting the residual against \hat{y}_i , so that is all for today.

Thank you.