

Regression Analysis and Forecasting
Prof. Shalabh
Department of Mathematics and Statistics
Indian Institute of Technology-Kanpur

Lecture 19

Software Implementation of Multiple Linear Regression Model using MINITAB

Welcome to the lecture you may recall that in the earlier lectures we had considered the set up of multiple linear regression model and we discussed different types of aspects related to the fitting of model judging the goodness of fit of a model and finally the objective was to get a good fitted model. Now let us take a simple example to understand whatever we have done in the earlier lectures, right. So if you see here am going to consider here an example where I am considering 27 observations.

(Refer Slide Time: 00:57)

Example:

Drinking water demand (in thousands kilo litre) depends on temperature and humidity level

y : Drinking water demand

X_1 : temperature

X_2 : humidity level

Model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i, i = 1, 2, \dots, 27$$

Drinking water demand (in thousands lit)	Temperature (in celsius)	Humidity Level (in %)
33710	23	16
31666	25	15
33495	25	19
32758	26	20
34067	27	20
36069	28	25
37497	30	25
33044	26	29
35216	29	28
35383	32	30
37066	33	30
38037	34	33
38495	35	35
39895	38	39
41311	39	40
42849	42	50
43036	43	55
43873	44	60
43923	45	70
45078	45.5	75
46935	45	80
47951	46	85
48085	44	94
48003	44	90
48080	41	92
42924	37	94
48061	40	95

And these observation have been obtained on the demand of drinking water, which is measured in 1000 kilo liter and we believe that this demand depends on the temperature that is the weather temperature and humidity level. You may also recall that in the case of linear regression model we also considered a similar example in which I have taken only variable that is temperature.

And now I am going to consider the second variable as humidity level and that is measured in percentage. So the model in which we are interested now here is $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$ and we had got 27 observation so I can express the model as $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i$.

(Refer Slide Time: 02:08)

$y = X\beta + \epsilon$
 $n \times 1$ $n \times k$ $k \times 1$ $n \times 1$
 $n = 27, k = 3$ including intercept term

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i,$$

$$i = 1, 2, \dots, 27$$

$$n = 27, k = 3$$

33710	1	23	16
33666	1	25	15
33495	1	25	19
32758	1	26	20
34067	1	27	20
36009	1	28	25
37497	1	30	25
33044	1	26	29
35216	1	29	28
35583	1	32	30
37066	1	33	30
38037	1	34	33
38495	1	35	35
38895	1	38	39
41311	1	39	40
42849	1	42	50
43038	1	43	55
43873	1	44	60
43923	1	45	70
45078	1	45.5	75
46935	1	45	80
47931	1	46	85
46085	1	44	94
48003	1	44	90
45050	1	41	92
42924	1	37	94
46061	1	40	93
47093	1	38	97
48517	1	35	96
49288	1	36	97

$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}$

 β_0 - int. term
 β_1, β_2 - slope param.

Now first let me explain how these things can be expressed in the framework of a multiple linear regression model. So here the model which we had considered was expressed as $y = X\beta + \epsilon$ where y was a n cross one vector of observations on response variables and X was a matrix of order n by k of k independent variable having n observation on each of the variable and β was a k cross one vector of regression coefficients and ϵ was a n cross one vector of random error components.

So now in this case you can see that here $n=27$ and k is equal to here actually three including the intercept term, okay, so now I can express all the 27 observations on y here like this which is a 27 cross one vector and there are two variables x_1 and x_2 on which we have got 27 observations and all these observations including the intercept term can be expressed as a matrix X which is of order 27 by 3.

And similarly β contains three elements first is intercept term and other two are slope parameters, β_1 is the slope parameters associated with x_1 and β_2 is the slope parameter associated with x_2 . Now let us try to use software which can analyze this data, so as discussed earlier we had seen that there are different types of statistical software available and some are paid software and some are free software among the paid software Minitab, Systat, SPSS, SAS and all such softwares are quite popular.

Among the free software R or GRETL these are two popular softwares among others, right. So in this lecture we are going to use the software Minitab, but definitely if you use any other software that will also give you the same outcome there can be some difference in the format that some software may give ANOVA table first and then the regression coefficient and some software may give the regression coefficient first and then the ANOVA table and so on

So there are only minor difference among the output and outcome of software, so that doesn't make any difference. Now before we go forward let me explain you what is our objective, my objective is not to find a model here by objective is that first we to expose this dataset to software and whatever is the outcome of software we have to understand that thing. So the first question is this, what different numerical values are indicating and how they have been computed and they correspond to wish of the quantity that we have covered in the lecture.

And after that we will try to analyze all the things one by one and finally we will try to combine all our outcomes together and then we will have a final fitted multiple linear regression model, right. Let us now start with this thing and I already have entered the data in this Minitab software.

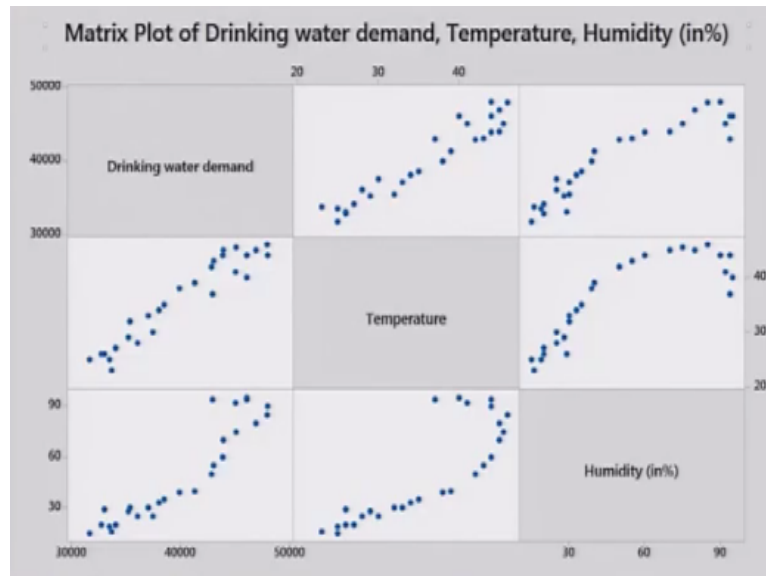
(Refer Slide Time: 06:00)

	C1	C2	C3	C4	C5	C6	C7
	drinking water demand	Temperature	Humidity (%)				
1	11710	23.0	58				
2	11660	25.0	53				
3	11440	25.0	57				
4	12750	26.0	30				
5	14507	27.0	35				
6	16000	28.0	25				
7	17407	30.0	25				
8	17044	26.0	29				
9	17210	26.0	28				
10	15000	32.0	30				
11	17060	33.0	30				
12	18007	34.0	13				
13	18405	35.0	25				
14	18805	36.0	38				
15	41111	39.0	40				
16	42009	42.0	50				

So you can see here that in the first column I have entered the data of drinking water demand the second column has the temperature and third column has the humidity. Now first of all we have to ensure that can we really use here multiple linear regression model, so we try to plot them since this is a multiple linear model so we would like to go for multiply scatter diagram or in the case of Minitab this is called as matrix plot.

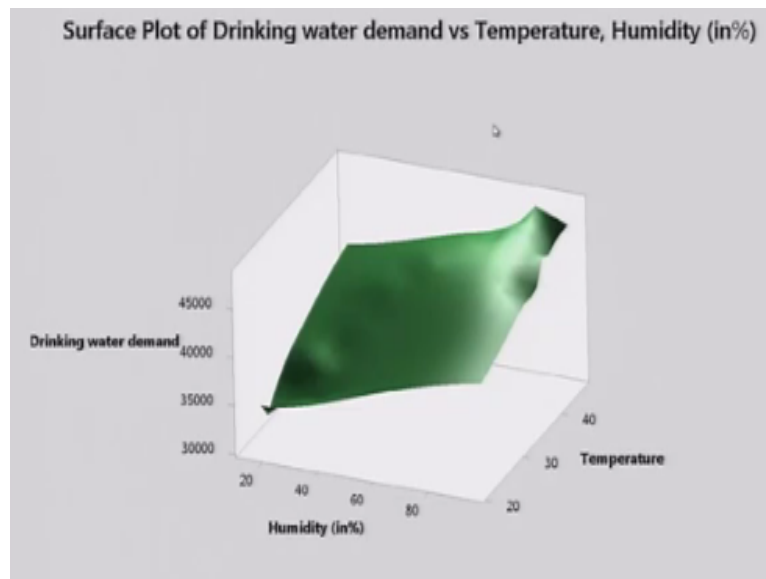
Right, so now there are different options available for the matrix plot and then we are just going to opt for a simple matrix plot and we can see that here I have chosen here all the three variables c1, c2, c3 and once I try to plot it I get here a scatter diagram like this one. Okay, let us try to first live it a side and try to make a some other things and at the end I will try to copy all the commands and then we will try to discuss them one by one.

(Refer Slide Time: 06:52)



Now since we have here only two variables x_1 and x_2 and third variable is here y , so one option we have that we can also make a three dimensional plot and so we try to go for a three dimensional plot, which is based on the wire. So here I try to using the Z axis is try to use the drinking water variable on y temperature and on X axis the humidity and we get here this type of three D surface.

(Refer Slide Time: 07:42)



Now let us try to store it and we will try to analyze it later, and similarly there is another option here that we can have here a three D surface plots, if try get the surface plots here that will give us this type of surface. So I am just trying to prepare the ground and we would like to be confident that once I fit a linear regression model we should be confident, yes this is situation where a multiple linear regression can be fitted.

Obviously when have here more than three variables all together than the option of three dimensional surface plot is not available and then in that case we have different options that we can take say three variables at a timed and then we try to combine all the individual inferences together and finally the objective is to take one regression in terms of obvious or no that can we fit here a multiple linear regression model or not.

Right, so these figures are giving us a sort of confidence that, yes a multiple linear regression model can be fitted over here. So now let me try to fit here a multiple linear regression model so you can see that in Minitab it is not difficult and it is a menu driven software. So now here I am taking my response variable as a drinking water demand and I have taken here two independent variable as temperature and humidity and here we do naught have any categorical variable which takes the values in terms of say zero one or something like that.

Let me try to explain what are the different meanings which are available to us, so here you can see that there is possibility here that one can also at the interaction term in multiple linear regression model once we decide yes two variable are interacting then the observation on

interaction are difficult to find in practice so what we try to do we simply try to multiply the values of the two corresponding variable and we obtain the observations on interaction.

But here in this case we are not considering the interaction term so we ignore it. Now if you see here we have here different options now I am setting here my this confidence level for all intervals as 95 that means the level of significant alpha is here 5% or point 05 now when am going for confidence interval I can have one sided interval are two sided interval so we decide in the case of in the of linear regression analysis we are usually interested in a two sided confidence intervals, so I give this option over here.

Well, now when we want o find out the sum a square there are two options that I can find out the adjusted sum of a square or say sequential sum of squares, but in the case of linear regression analysis what we have considered we are interested in the simple sum of a square are that is called here as a adjusted sum of a squares, right adjusted had some other meaning but that is beyond the scope of this lecture.

Another option here is this Box-Cox transformation that we have not considered in this lecture, but this is a sort of transformation that is applied to the independent variables if y is not really related to x, but it is linearly related to some other power of x something like square root of x or x square something like this, but we are not considering here so we have just using to define the level of alpha and the type confidence interval and the type of sum of squares.

Since we are not using here the categorical variables, so coding is not important and we are not using here the stepwise regression, stepwise regression is another type of regression, which we have not described in this course. Now there is another option to have here the graphs, so you can see here that there are different options, histogram of residues, normal probability plot, residual versus plot, residual versus order.

And you may recall that when we discuss the diagnostic for our multiple linear regression model we had considered all the things but in this case I would be more interested in knowing whether the observations are originating form a normal distribution or not and then I would try to have a plot for the residual versus fits, so that I can decide for the presents of heteroskedasticity in my model.

Similarly for the residual versus order we have a time series plot so we can have all these things, now here if you see now we have here three options that we can make this residual plots with respect to simple residuals that we have defined by ϵ_i are the standardized residual that we have define say as d_i or say r_i say studentized residual or standardize residuals and deleted residual that we had defined as e_i and i inside the bracket.

So here we would like to have the regular one, and then we would also like to see that about what happens when we are going for a standardized one, ideally there should not be any difference. So in the first shot I will try to obtain the these plots for the regular residuals and then I will try to repeat this analysis and I will obtain the plots for the standardized residuals and there is another option that I can have all the four graphics in one but anyway that doesn't make any difference.

Now the another option is here results that what type of results I would like to have, one thing is the method, yes we would like to know that what type of method we have used like least square principal or something like this then we would also like to know what are the outcomes of analysis of variance and we would also like to know the model summary that consist of fourth the statistics like r^2 adjusted r^2 PRESS and so on and we would also like to store the values of a coefficients that is the values of β_0 , β_1 and β_2 , and we would also like to have a regression equation that is obtained after fitting the model.

And we would also like to know the fits and diagnostic if you remember we had done different types residuals and different types of statistic like as h_{ii} is the i th diagonal element of hat matrix, deleted residuals DFFITS and so on. Now I have here two options that whether that I would like to fits and diagnostic for all the observation and or say only for those observation which are found to be unusual.

But since here we are learning, so we would like to see what happens to all the observation and then there is a statistics what is here called as a Durbin-Watson statistics we have not done this thing in detail but still I would try to its outcome and then I will try to explain you how to interpret it. Durbin-Watson test is a test to check the presence of first order

autocorrelation in the data that means increase if the data are affected by the time then we have time series analysis.

And in the case of time series analysis we try to test whether the observations have first order autocorrelation or something else. Right, now in case if you want you can also store these results whatever you have obtained here, right so this can be stored in the worksheet that is opened here and one may use it. So we would try to store to just illustrate that how this actually works.

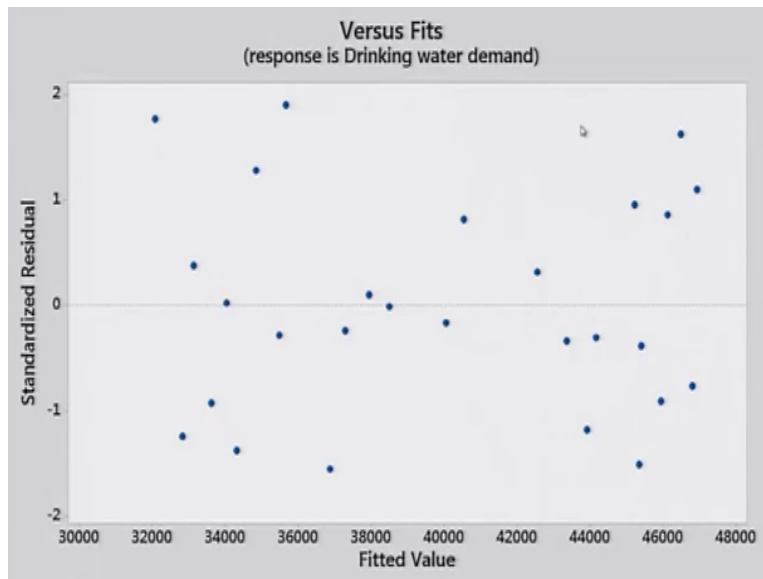
So now once I go for this analysis let us try to see what happens, so now you can see here that we had obtained here a normal probability plot so we try to store it and we have got here a plot between residuals and fitted value, okay. So let us try store it, now if you observe over here this is the outcome here we have analysis of various results in here and here we have the details about the model summary which is consisting about r^2 , adjusted r^2 , PRESS and r^2 prediction.

And similarly here we are getting the value of the regression coefficient their confidence interval various inflation factor values for the t statistics and here we are getting the fitted regression equation, and here there are different values of diagnostics and since we had given that we would like to have this diagnostic for all the observation so this is giving here say the fit standard error of it confidence interval for the \hat{y} residual, standardized residual, deleted residual h_i and Cook's D statistics here you can see.

Right, and this is available for all the 27 observations and similarly the DFFITS statistics is also computed here for all the observation and at end we have here Durbin-Watson statistics and we try store these results separately and then we try to analyze it, but before go for the further analysis let us also try to obtain the this residual plots for the standardized case.

So now I can choose here this option standardized and let me do this analysis, so you can see here that we had now here these two graphics when the residuals are standardized this is the normal probability plot and this here the standardized residual versus fitted value plot.

(Refer Slide Time: 18:58)



(Refer Slide Time: 19:04)

$y = X\beta + \varepsilon$
 $n \times 1 \quad 70 \times 2 \quad 27 \times 1 \quad n \times 1$
 $n = 27, k = 3$ including interest term

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i,$$

$$i = 1, 2, \dots, 27$$

$$n = 27, k = 3$$

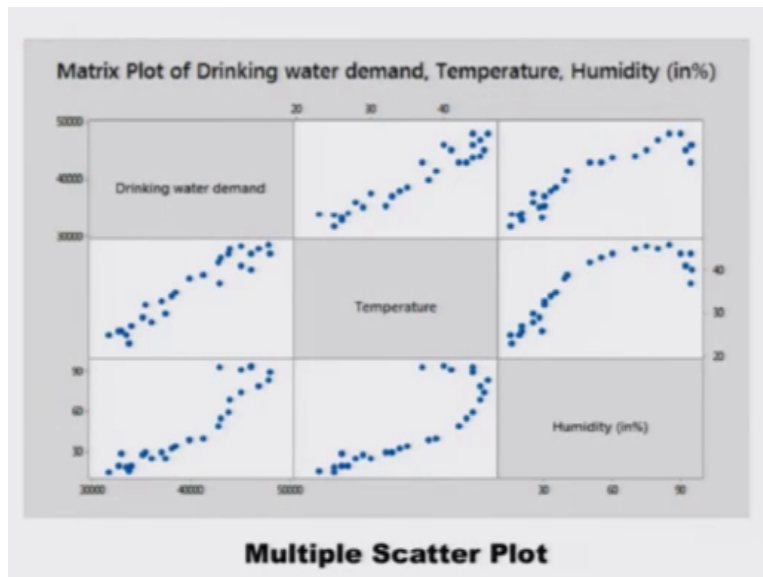
33710	1	23	16
31666	1	25	15
33495	1	25	19
32758	1	26	20
34067	1	27	20
36069	1	28	25
37497	1	30	25
33044	1	26	29
35216	1	29	28
35383	1	32	30
37066	1	33	30
38037	1	34	33
38495	1	35	35
38895	1	38	39
41311	1	39	40
42849	1	42	50
43038	1	43	55
43873	1	44	60
43923	1	45	70
45078	1	45.5	75
46935	1	45	80
47951	1	46	85
46085	1	44	94
48003	1	44	90
45050	1	41	92
42924	1	37	94
46061	1	40	95
47093	1	38	97
48517	1	35	96
49288	1	36	97

$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}$

β_0 - int. term
 β_1, β_2 - slope para.

So now we have stored all this information over here and we would try to understand the interpretation and other things for these results.

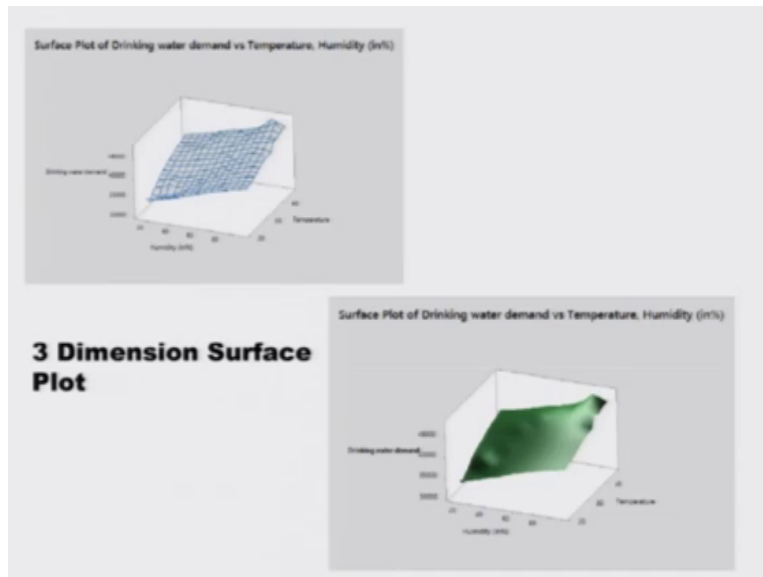
(Refer Slide Time: 19:15)



So you see this is the matrix plot of the three variables y , x_1 and x_2 . One can see here this is a plot between y and x_1 and one can see here that it is showing as a linear trend so there is no issue similarly this is a plot between y and x_2 and this is also showing here a linear trend and similarly this is also a plot between x_1 and y this is showing a linear trend and similarly this is also a plot between x_2 and y this is showing us a linear trend.

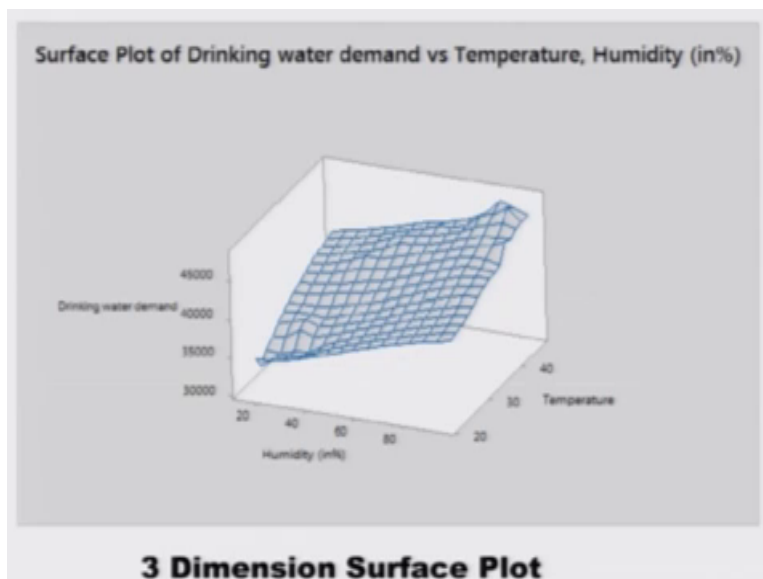
Now if you try see here this is a plot between x_1 and x_2 , you can see here ideally we assume that x_1 and x_2 should be independent so we should have a scatter diagram like this one, but here you could see it is trying to show here a sort of linear relationship. So let us see whether this is reflected in our regression analysis or not, but more or less this is giving us a surety that at least the relationship between y and x_1 is linear, relationship between y and x_2 is linear. So we can believe that the relationship between y and x_1 jointly will also be linear.

(Refer Slide Time: 20:43)



Now these are the surface plots that we had obtained you can see here that in deep case this is also showing us sort of increasing a trend and this is the same story her also this also showing a sort o f increasing trend so these three dimensional plot which are useful particular in this example are also giving us confidence that a multiple linear regression model can be fitted in this case, right okay so this is the same thing that we had obtained it is in a bigger shape.

(Refer Slide Time: 21:08)



(Refer Slide Time: 21:26)

$$\hat{\beta} = (X'X)^{-1}X'y = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}$$

$$\widehat{Cov}(\hat{\beta}) = \begin{pmatrix} \widehat{Var}(\hat{\beta}_0) & \widehat{Cov}(\hat{\beta}_0, \hat{\beta}_1) & \widehat{Cov}(\hat{\beta}_0, \hat{\beta}_2) \\ & \widehat{Var}(\hat{\beta}_1) & \widehat{Cov}(\hat{\beta}_1, \hat{\beta}_2) \\ & & \widehat{Var}(\hat{\beta}_2) \end{pmatrix} = \hat{\sigma}^2 (X'X)^{-1} = \hat{\sigma}^2 C$$

$$\widehat{Var}(\hat{\beta}_j) = \hat{\sigma}^2 c_{jj} \quad j=0,1,2$$

Standard error of $\hat{\beta}_j = \sqrt{\widehat{Var}(\hat{\beta}_j)}$

Confidence intervals for $\hat{\beta}_j$

$$\hat{\beta}_j \pm t_{\frac{\alpha}{2}, n-k} se(\hat{\beta}_j)$$

Test of hypothesis $H_0: \beta_j = 0$

$$T_j = \frac{\hat{\beta}_j - 0}{se(\hat{\beta}_j)} \sim t(n-k)$$

Now after this we are going to discuss about the outcomes of regression analysis, but before that let me briefly describe some results so that it is easier for us to understand you may recall that we had estimated beta hat say x transpose x whole inverse x transpose y .

(Refer Slide Time: 21:30)

$$\hat{\beta} = (X'X)^{-1}X'y = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}$$

$$\widehat{Cov}(\hat{\beta}) = \begin{pmatrix} \widehat{Var}(\hat{\beta}_0) & \widehat{Cov}(\hat{\beta}_0, \hat{\beta}_1) & \widehat{Cov}(\hat{\beta}_0, \hat{\beta}_2) \\ & \widehat{Var}(\hat{\beta}_1) & \widehat{Cov}(\hat{\beta}_1, \hat{\beta}_2) \\ & & \widehat{Var}(\hat{\beta}_2) \end{pmatrix} = \hat{\sigma}^2 (X'X)^{-1} = \hat{\sigma}^2 C$$

$$\widehat{Var}(\hat{\beta}_j) = \hat{\sigma}^2 c_{jj} \quad j=0,1,2$$

Standard error of $\hat{\beta}_j = \sqrt{\widehat{Var}(\hat{\beta}_j)}$

Confidence intervals for $\hat{\beta}_j$

$$\hat{\beta}_j \pm t_{\frac{\alpha}{2}, n-k} se(\hat{\beta}_j)$$

Test of hypothesis $H_0: \beta_j = 0$

$$T_j = \frac{\hat{\beta}_j - 0}{se(\hat{\beta}_j)} \sim t(n-k)$$

So in this case this is going to be like beta0 hat beta1 hat and beta2 hat, and a three cross one vector of these three values and similarly we have obtained the covariance matrix of beta hat, and which in this case will variance of beta0 beta naught variance of beta1 hat and variance of beta2 hat, and on the off-diagonal elements this will be covariance between beta0 hat beta1 hat and covariance between beta0 hat and beta2 hat.

And here the covariance between beta1 hat and beta2 hat so this matrix will look like this and but our interest here is to go for its estimate, so the estimates and this value was obtained as

sigma square x transpose x whole inverse but here we are interested in the estimate of covariance matrix, so that we had obtained like this so this will be the estimated values and this will sigma square hat x transpose x.

Okay, and you may recall that we had denoted for example this sigma square some matrix here see c where we see that variance of beta j and its estimator is given by sigma square hat cjj, j goes around here in case of zero1 and 2. So that is the symbol that we are going to use here and when we are trying to find out the standard error then standard error is of beta hat j this is simply the positive e square root of estimate of variance of beta hat j.

Then we also had obtained the confidence interval, the confidence interval were obtained as say for say beta hat j this was obtained as beta hat j+- t alpha by 2 and - k and standard error of beta hat j and for the test of hypothesis when we are interested in testing the hypothesis h naught beta j=0 then the t statistics for the jth regression coefficient was obtained as beta h j -0 divided by standard error of beta hat j, and this was following a t distribution with n-00 k degrees of freedom.

(Refer Slide Time: 24:39)

Coefficients	Coef	SE Coef	95% CI	T-Value	P-Value	VIF
Constant	$\hat{\beta}_0 = 21361$	1206 $SE(\hat{\beta}_0)$	(18871, 23851) $\hat{\beta}_0 \pm t_{\alpha/2, 24} SE(\hat{\beta}_0)$	17.71	0.000	$VIF = \frac{1}{1-R_j^2}$ Inflation factor LCVIF < 3 ⇒ moderate collinearity
Temperature	$\hat{\beta}_1 = 412.3$	46.6 $SE(\hat{\beta}_1)$	(316.2, 508.5) $\hat{\beta}_1 \pm t_{\alpha/2, 24} SE(\hat{\beta}_1)$	8.85	0.000	
Humidity (in%)	$\hat{\beta}_2 = 77.5$	12.6 $SE(\hat{\beta}_2)$	(51.6, 103.5) $\hat{\beta}_2 \pm t_{\alpha/2, 24} SE(\hat{\beta}_2)$	6.18	0.000	

Regression Equation
 Drinking water demand = 21361 + 412.3 Temperature + 77.5 Humidity (in%)
 $y = 21361 + 412.3X_1 + 77.5X_2$

Model Summary

S	R-sq	R-sq(adj)	PRESS	R-sq(pred)
994.608	96.73%	96.46%	30407053	95.81%

Handwritten notes: $H_0: \beta_j = 0$, $T_j = \frac{\hat{\beta}_j - 0}{SE(\hat{\beta}_j)}$, H_0 is rejected if $P < \alpha$, $\alpha = 0.05$.

Now you will see that the values which we have obtained earlier analytically now they are being computed easily by the software so let us try to understand here the this outcome I have divided the entire outcome into separate sheets for that we can understand them clearly. Let us go one by one, so when we are trying to see here this coefficient, coefficient is nothing but your beta hat vector so this is here beta hat zero this is the constant.

Constant means the estimator of intercept term and the temperature means here the regression coefficient corresponding to the independent variable temperature x_1 and similarly this is here the value of $\hat{\beta}_2$, which is associated with the second independent variable humidity. Now this is the standard error, so these values are giving us the idea of standard error of different coefficient, so this is standard error of $\hat{\beta}_0$, this the standard error $\hat{\beta}_1$ and this is the standard error $\hat{\beta}_2$.

The next column is about the 95% confidence interval, so this is the lower confidence limit and this is the upper confidence limit for $\hat{\beta}_0$ this is provided by this expression $t_{\alpha/2}$ with degrees of freedom $n - k$, which is 27-324 and standard error of $\hat{\beta}_0$ and similarly this is also a confidence interval for $\hat{\beta}_1$ obtained like where $\alpha = 0.05$.

And this is the confidence interval for $\hat{\beta}_2$ this is obtained at $\alpha = 0.05$ with 24 degrees of freedom and standard error $\hat{\beta}_2$. So these things give us this idea, so before we try to go for the next column let us try to understand what is the interpretation, the interpretation is that on the basis of given set of data we can say that the value of intercept term here is 21361 this is the point estimate.

And the confidence interval for $\hat{\beta}_0$ is like this, this means the value of β_0 is going to lie between 18871 and 23851, so you can see that the point estimate is lying within the confidence interval and similarly the value of this $\hat{\beta}_1$ is indicating that the point estimate is 412.3 where as the confidence interval is 316.22508.5.

So one can also see here that this is going to lie between the confidence interval and the same thing is for this $\hat{\beta}_2$ also, right. The physical interpretation of $\hat{\beta}_0$ is the following; that when $x_1=0$ and x_2 is also $=0$ then the average value of y is going to be 21361 right, and similarly the physical interpretation of $\hat{\beta}_1$ is that when there is a unit change in the value of x_1 .

For example in this case x_1 is that temperature, so when that the temperature changes by 1 degree centigrade then the change in the average value of y is 412.3 that means once the temperature changes by 1 degree centigrade then the demand for water on an average increases

by 412.3000 kilo liter the question is why I am saying it increasing because one can see here that the sign of beta one hat is positive.

In case if the sign would had been negative then I would say that the demand would decrease by 412.3000 kilo liter when the temperature changes by one degree. Okay, similarly here also in the case of humidity also it is trying to say that when there is a change of 1% in the humidity it increases the average value of y changes by 77.5 units.

That is when there is a 1% change in the humidity on an average the demand of drinking water in cases by 77.5000 kilo liter, right and these are the indicators like as standard errors, this is giving us an idea about the spread of beta0 hat beta1 hat and beta2 hat, okay. Now we come to this column t value, t value is about the test of hypothesis that $H_0: \beta_j = 0$.

So for example this and here $t_j = \hat{\beta}_j - 0$ divided by standard error of beta hat j, so this value is beta hat naught divided by standard error of beta hat naught this is the value of t statistics for the intercept term and similarly this is also the value of beta1 hat divided by standard error of beta1 hat that we had denoted by here as t one and similarly this is also a value which is beta hat 2-0 divided by standard error of beta2 hat, so these are the values of t statistics.

Now as I told you earlier that in order to test the hypothesis we have two options either we can obtain the critical values or the tabulated values from the probabilities of t tables or in software there is another option to use the p values and we had learned that the result is that H_0 is rejected if p value is less than alpha that is the level of significance so here we have taken $\alpha = 0.05$.

So you can see here that the p value in all the three cases they are very, very close to 0, so one can say that in this case all the three null hypothesis that is $H_0: \beta_0 = 0$, $H_0: \beta_1 = 0$ and $H_0: \beta_2 = 0$ all are rejected. So that means none of the values beta0, beta1 or beta2 is 0 or close to 0. This means the variables which we have considered as a temperature and humidity.

They are contributing in the model and they are helping in explaining the variation in the value of y. The last column is various inflation factor you may recall that we had discussed this diagnostic when we wanted to test the assumption of rank of $X = k$ that means all x_1, x_2, \dots, x_k are independent or not and we had defined that $VIF = \frac{1}{1 - r_j^2}$.

Where r_j^2 is the coefficient of determination that is obtained by regressing that j th explanatory variable x_j over the remaining $k-1$ explanatory variables, and if you remember we had given the criteria that in case if VIF is lying between one and three this will be indicating that the independent variables are moderately correlated and if it is more than three one can consider that they are highly correlated.

(Refer Slide Time: 33:24)

$$\hat{\beta} = (X'X)^{-1}X'y = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}$$

$$\widehat{Cov}(\hat{\beta}) = \begin{pmatrix} \widehat{Var}(\hat{\beta}_0) & \widehat{Cov}(\hat{\beta}_0, \hat{\beta}_1) & \widehat{Cov}(\hat{\beta}_0, \hat{\beta}_2) \\ \widehat{Cov}(\hat{\beta}_1, \hat{\beta}_0) & \widehat{Var}(\hat{\beta}_1) & \widehat{Cov}(\hat{\beta}_1, \hat{\beta}_2) \\ \widehat{Cov}(\hat{\beta}_2, \hat{\beta}_0) & \widehat{Cov}(\hat{\beta}_2, \hat{\beta}_1) & \widehat{Var}(\hat{\beta}_2) \end{pmatrix} = \hat{\sigma}^2 (X'X)^{-1} = \hat{\sigma}^2 C$$

$$\widehat{Var}(\hat{\beta}_j) = \hat{\sigma}^2 c_{jj} \quad j = 0, 1, 2$$

Standard error of $\hat{\beta}_j = \sqrt{\widehat{Var}(\hat{\beta}_j)}$

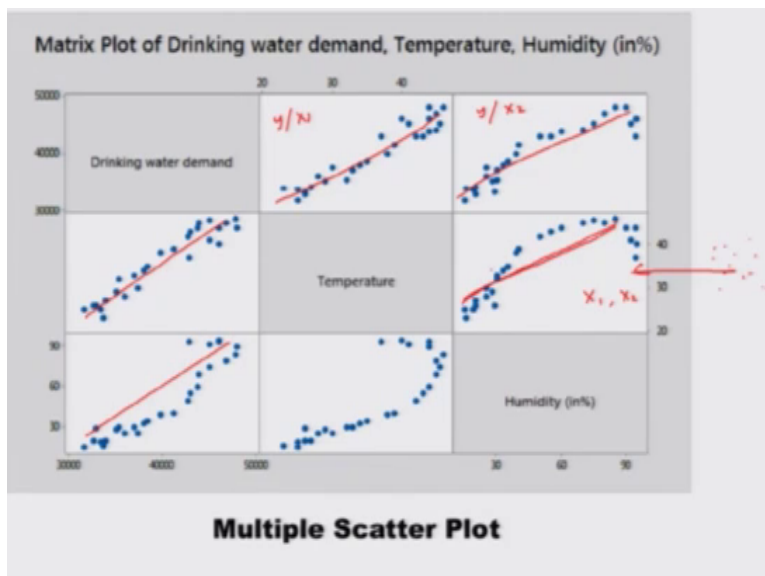
Confidence intervals for $\hat{\beta}_j$

$$\hat{\beta}_j \pm t_{\frac{\alpha}{2}, n-k} se(\hat{\beta}_j)$$

Test of hypothesis $H_0: \beta_j = 0$

$$T_j = \frac{\hat{\beta}_j - 0}{se(\hat{\beta}_j)} \sim t_{(n-k)}$$

(Refer Slide Time: 33:43)



Now if you see here in this multiple scatter diagram which I indicated her you can see that that there is linear trend here a sort of linear trend so one can see here that x_1 and x_2 are not really uncorrelated, but they have got certain degree of linear relationship and this is being indicated in this software here in this value 3.041 like this right okay. So you can see that our statistical analysis is also capturing this aspect.

(Refer Slide Time: 34:02)

Term	Coef	SE Coef	95% CI	T-Value	P-Value	VIF
Constant	$\hat{\beta}_0 = 21361$	1206	(18871, 23851)	17.71	0.000	
Temperature	$\hat{\beta}_1 = 412.3$	46.6	(316.2, 508.5)	8.85	0.000	3.41
Humidity (in%)	$\hat{\beta}_2 = 77.5$	12.6	(51.6, 103.5)	6.18	0.000	3.41

Regression Equation
 $y = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2$
 Drinking water demand = 21361 + 412.3 Temperature + 77.5 Humidity (in%)

Model Summary

R-sq	R-sq(adj)	PRESS	R-sq(pred)
99.4608	96.46%	30407053	95.81%

Handwritten notes: $H_0: \beta_j = 0$, $T_j = \frac{\hat{\beta}_j - 0}{SE(\hat{\beta}_j)}$, H_0 is rejected if $p < \alpha$, $\alpha = 0.05$. VIF = $\frac{1}{1 - R_j^2}$. Variance Inflation factor. $CVIF < 3 \Rightarrow$ no multicollinearity. $\hat{\beta}_j \pm t_{\alpha/2, n-k} SE(\hat{\beta}_j)$. $R^2 = 1 - \frac{SS_{res}}{SS_T}$. $R^2 = 1 - \frac{(n-1)MS_{res}}{(n-k)MS_{reg}}$. $R^2_{adj} = 1 - \frac{PRESS}{SST}$. $\hat{\sigma}^2 = \frac{SS_{res}}{n-k}$. $E(\hat{\beta}_j) = \beta_j$ (Assuming no multicollinearity).

Now based on this regression coefficients outcome we can now formulate our fitted regression model as $y = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$. Now after this we have here some more model summary this value here as S this is nothing but your square root of sigma square hat, and sigma square hat was obtain by SS residual divided by degrees of freedom n-k.

(Refer Slide Time: 37:31)

Source	DF	Seq SS	Contribution	Adj SS	Adj MS	F-Value	P-Value
Regression	2	702690691	96.73%	702690691	351345346	355.17	0.000
Temperature	1	664968717	91.54%	77519858	77519858	78.36	0.000
Humidity (in%)	1	37721974	5.19%	37721974	37721974	38.13	0.000
Error	24	23741875	3.27%	23741875	989245		
Total	26	726432566	100.00%				

Handwritten notes: SS : Sum of Squares, MS : Mean Squares, F : F-statistic. $SS_{reg} = SS_{reg(x_1)} + SS_{reg(x_2)}$. $MS_{reg} = \frac{SS_{reg}}{DF}$. $F = \frac{MS_{reg}}{MS_{res}}$. SS_T : Total SS. $H_0: \beta_1 = \beta_2 = 0$, $\alpha = 0.05$, $P\text{-value} = 0.00 < \alpha \Rightarrow H_0$ is rejected.

This is our here R square R square was a statistics that was used to test the goodness of fit and this we had defined as $1 - \text{SS residual} / \text{total sum of squares SST}$, so one can see here that this is 96.73, so one can be happier that the model which we have fitted is reasonably good and essentially the variables x_1 and x_2 are capable of explaining say at least 96.73% variation in the values of y .

Next quantity is here R square adjusted you may recall that we had define this quantity like R^2_{adj} and this was define as $1 - \frac{1}{n-1} \frac{\text{SS residual}}{\text{SST}}$. So this is again say 96.46 which is again indicating that our model is nearly 96% good and you will always see that R square is always greater than the value of R square adjusted, but usually in practice this difference is very small.

Next here is PRESS, the PRESS was the statistics that was given by $\sum_{i=1}^n \epsilon_i^2$. You may recall that these ϵ_i hats were PRESS residuals that means the different between y_i and say \hat{y}_i which is obtain after deleting the i th observation and based on this PRESS statistics we had defined the R square prediction and which is obtain here by $1 - \text{PRESS} / \text{total sum of square SST}$.

And this R square predication was the quantity that was helping us in deciding whether the fitted model can be used for prediction or not and if used how it is going to perform. So this value here 95.81%, so it is simply indicating that this model is nearly 95% good for predication. So now let us try to concentrate on the on other part which is analysis of variance.

So here degrees of freedom means is denoted by here df , these are degrees freedom and here the software is computing other component sequential sum of squares and it is also computing the contribution of sequential sum of square, which we have not covered in our lecture, but I will try to explain you its importance over here, what we considered in our lecture is that we have obtained the sum of squares which is denoted here has adjusted of square.

So this is essentially in our terminology this is SS sum of a squares and this is in our terminology this is mean squares and this is value of F- statistics and this is the P-Value, now you can recall that we had divided the total sum of a squares in the case of analysis of

variance in two to components sum of square due to regression and sum of a squares due to error or sum of the square due to residual.

So this quantity here corresponding to this, this is giving us the value of SS res, that is the some of the square due to regression, and this is the quantity here which is trying to give us the value of SS res that is residual sum of squares. Now the software is also dividing this some of the due to regression into two components which is some of the squares due to regression due to the variable x_1 , and this is the sum the square regression due to the variable x_2 .

And yeah, and here also you can see this is the sequential sum of the square which they again divided it into two components due to x_1 due to x_2 and they are using it in understanding the contribution of x_1 and x_2 in the regression sum of the square, so here you can see this value is 91.54% that means out of the two variable temperature and humidity, temperature is contributing nearly 91.54% and humidity is contributing only 5.19%.

Right and the component of error here is only 3.27%, right, the sum of squares due to residual will remain the same as SS res, right. So if you try to see even at these three values the make sense whenever the temperature is high usually the demand of water will increase, but if the humidity is also high the demand of water also increases, but the rate of change is low.

So that make sense that the model which we are fitted and the values which we have obtained here they are making some sense. Now this value here is mean square due to regression, and this value here is, this is mean square due to residuals okay, and the same thing here they have obtain the mean square due to regression due to x_1 and here they have obtained it mean square due to regression due to the second variable x_2 .

And this value is the value of F-statistics that was obtain as MS res divide by MS res, so corresponding to this thing the p value is coming out to be 0.00, okay and here this the value of SST, this the total some of squares, and here you can see the degrees of freedom also here it is two, which was $=k-1$ that is $3-1$, and this degrees of freedom was $n-1$, so this 27-and this was a $n- k$, so this is 24 degrees of freedom.

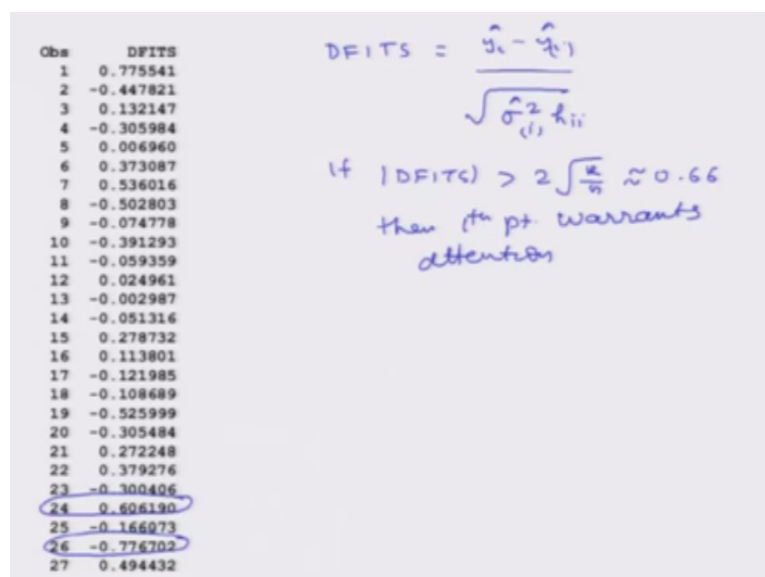
had denoted ϵ_i and this h_{ii} is the i th diagonal element of $h = X(X^T X)^{-1} X^T$, and this is our Cook's D statistics and so this in our notation we also have denote it has here d_i or d_i 's values are indicated by d_i .

Right, now we try to look at the interpretation of h_{ii} , so you may recall that we had concluded that if h_{ii} is greater than $2k/n$ which was a sort of guideline in our case this is 0.2 this would indicate that there are no leverage points. So you can see here in these values none of the values is very close to 2.2 expect here these values over here.

Right, and then we will see later on also that these two values are quite far away from the fitted regression line so this is also well captured by our defined statistics. Okay, now we come to the interpretation of d_i 's that is the Cook's D statistics. We had made a conclusion that that in case if d_i is greater than one then we would conclude that the data point is influential and it needs your attention.

But you can see here that none of the value here is close to one, so there is no question of having outlier and this was also indicated by our scatter diagram also that there are no such influential point which are trying to drag the fitted regression line towards them.

(Refer Slide Time: 47:14)

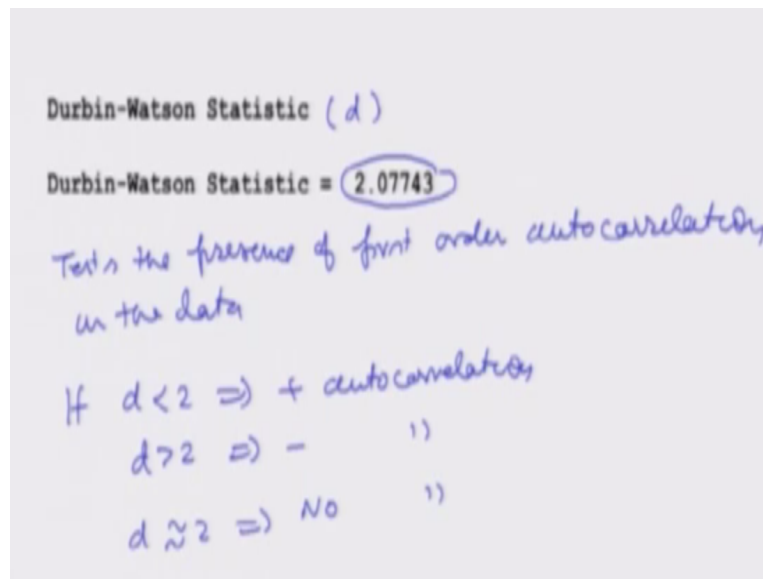


Right, okay, now we come to our next data, here it is trying to give us the value of the DFITS statistics and you might call that we had defined the DFITS statistics as y_i minus \hat{y}_i based on deleted residuals divided by σ_i^2 hat, based on deleted residuals h_{ii} and we had

conclude that if absolute value of DFITS is greater than twice of square root of k by n which in our case comes out to be 0.66 then ith point warrants attention.

So you can see here that here only there are two values here which are closed to or greater than point six, so here we have to be little bit careful that we will also see in the scatter diagram that two at the points at the end they are going towards downward direction, so this phenomenon is again correctly capture by our statistics and our regression analysis.

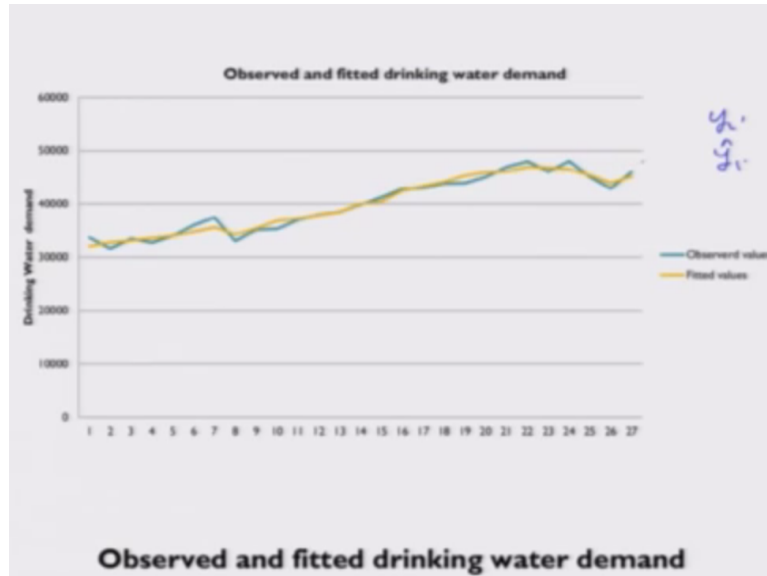
(Refer Slide Time: 48:32)



Right, now the next value is about Durbin-Watson statistics this is usually denoted by here d and Durbin-Watson statistics actually tests the presence of first order autocorrelation in the data and here the decision rule is very simple that if d is a smaller than two this indicates positive autocorrelation and if d is greater than two this indicates negative autocorrelation auto if d is close to two this indicates no auto correlation.

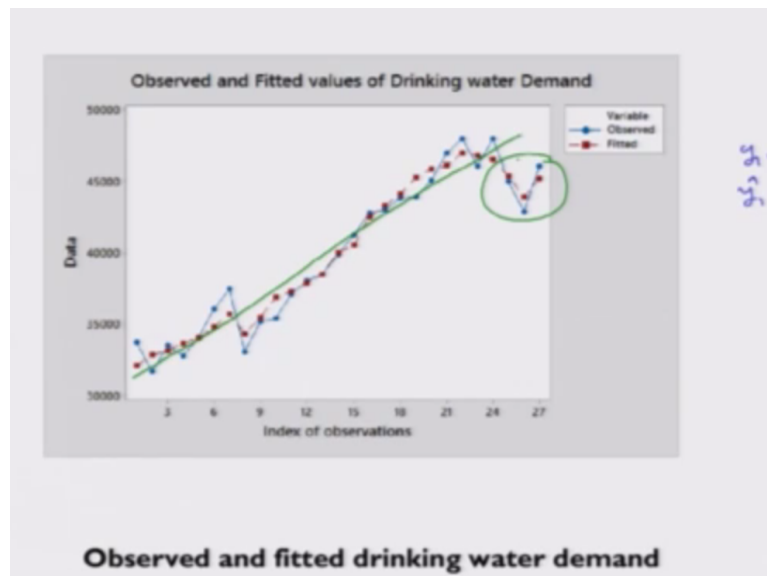
So now here in this case you can see that the value of Durbin-Watson statistics is close to two, so there is no autocorrelation of first order present in the data

(Refer Slide Time: 49:31)



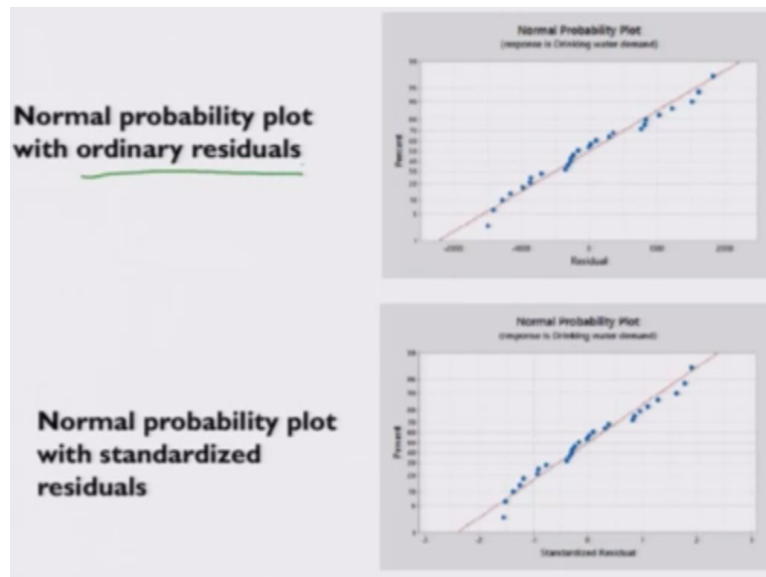
Right, now we try to analysis these values through our graphical procedures so here I have plotted the values of y_i and say \hat{y}_i , which I have obtained after fitting the model. One can see here that there is not much difference between the two values, so I can conclude that my model is well fitted and I even if I try to think it to use it for the prediction it is not difficult to say this model will continue further like this. Okay, and similarly in the next graphic also I have plotted these y_i and \hat{y}_i point by point.

(Refer Slide Time: 50:07)



So you can see her say from here to here there is a sort of linear trend and towards the end part over here this is showing a sort of trend towards the downward and that was actually captured in the d statistics and say residuals. Okay so this is also confirming about the properties of the data, right. Now this is the normal probability plot.

(Refer Slide Time: 50:44)



And first figure is based on the normal probability plots based on ordinary residuals and second figure this the normal probability plot using the standardized residual. So you may recall that we had concluded that if all such points which are lying over here as a blue dots if they are lying closed to the line nearly on a straight line then we can be sure that the observation are coming from a normal population.

And this phenomena is happening in both the cases whether we are using the ordinary residual or the standardized residual so this gives us an assurance that the assumption what we have made that this correct that our observation are originating from normal distribution.

(Refer Slide Time: 51:34)



Now these are the plot of residuals versus the fitted value, so here in the first case I have taken the ordinary residual and in the second case I have taken the standardized residuals, and

you can see that there is not much difference between the two. Right and from these two pictures we can conclude that all the points possibly they can be contained in a sort of horizontal band and there no pattern like outward opening funnel or inward opening funnel or there are no indication of having non linear relationship.

So from these two graphics we can conclude there is no problem of heteroscedasticity in the data and our variance remains constant. Okay now I have consider here an example and I have explained the complete statistical analysis whatever we have done in the earlier lectures in all details. Now please try to have a look at all these results, and try to understand yourself that which of the quantity is indicating what and how it is computed, once you know this thing then you will have more inside into data analysis.

Particularly in those situation when there is a problem because you are not going to deal in practice with 27 observation, but there can be several hundred observations and then these diagnostics these statistical results will help you in taking a decision that whether the fitted regression model is good or bad, you may not have an opportunity to look at the individual data points, and you would like to inspect the individual data points only in case of any problem.

So I stop here and in the next lecture I will try take one more example to illustrate some further aspects of multiple linear regression analysis, till then good bye.