

Business Statistics
Prof. M. K. Barua
Department of Management Studies
Indian Institute of Technology – Roorkee

Lecture – 58
Assumption of Regression

Hello friends. I welcome you all in this session. As you are aware, in previous session we were discussing about regression and correlation. We have seen that regression is basically association between two variables and correlation is nothing but it is a measure of strength of relationship between two variables. We did work out an example as well. In today's class, we are going to talk about assumptions which you should fulfill before applying regression to any data set.

(Refer Slide Time: 01:03)

The slide is titled "Assumptions of Regression L.I.N.E". The acronym "L.I.N.E" is circled in red. Below the title, there are four bullet points, each with a red arrow pointing to it and some underlines or red markings:

- Linearity
The relationship between X and Y is linear
- Independence of Errors
Error values are statistically independent
- Normality of Error
Error values are normally distributed for any given value of X
- Equal Variance (also called homoscedasticity)
The probability distribution of the errors has constant variance

Handwritten notes in red ink on the right side of the slide include: "observed - predicted", "F1 → 200", "190", "110", "210", and "110". There are also several red arrows and underlines throughout the slide.

At the bottom of the slide, there are logos for "IIT ROORKEE" and "NPTEL ONLINE CERTIFICATION COURSE".

So there are 4 basic assumptions, unless until you fulfill these assumptions you should not apply regression. Otherwise, you would be getting wrong results. In fact, you would get results without in fact looking for these assumptions but that could be wrong result right. So the first is linearity. So we are calling it as LINE, linearity independence of error term normality and equal variance right.

So first is there should be linear relationship between X and Y right. This is first and foremost assumption. If you do not have linearity of X and Y, if it is non-linear then you should convert it into linear function first and then apply regression. Independence of error terms so error values are statistically, they should be statistically independent. Error means it is

whatever is basically the difference between the observed value and actual value right, observed and actual right.

In fact, let me give an example of error term. Let us say if you forecast something, let us say the forecast of a particular product in let us say next month would be this much right and the actual demand turns out to be 90 right. So will say there is an error of 10 units right, is not it? So error 10. Now let us look at so this positive error let us say if your forecast is 100 units but actual sales comes out to be 110.

So the error would be -10 so error values are statistically independent right. So the error means what error it is a difference between actual and what you have predicted right. Normality of error right, so error values should normally be distributed for any given value of X, so we will see. In fact, we have seen several methods of checking normality of data; let us say normal probability plot, beat and normal distribution curve and so on right.

The fourth one is equal variance right, so the probability distribution of the error has constant variance. So these are 4 basic assumptions which you must fulfill whenever you solve any question using regression.

(Refer Slide Time: 03:56)

Residual Analysis

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

- The residual for observation i , e_i , is the difference between its observed and predicted value
- Check the assumptions of regression by examining the residuals
 - Examine for Linearity assumption
 - Evaluate Independence assumption
 - Evaluate Normal distribution assumption
 - Examine for constant variance for all levels of X (homoscedasticity)
- Graphical Analysis of Residuals
 - Can plot residuals vs. X

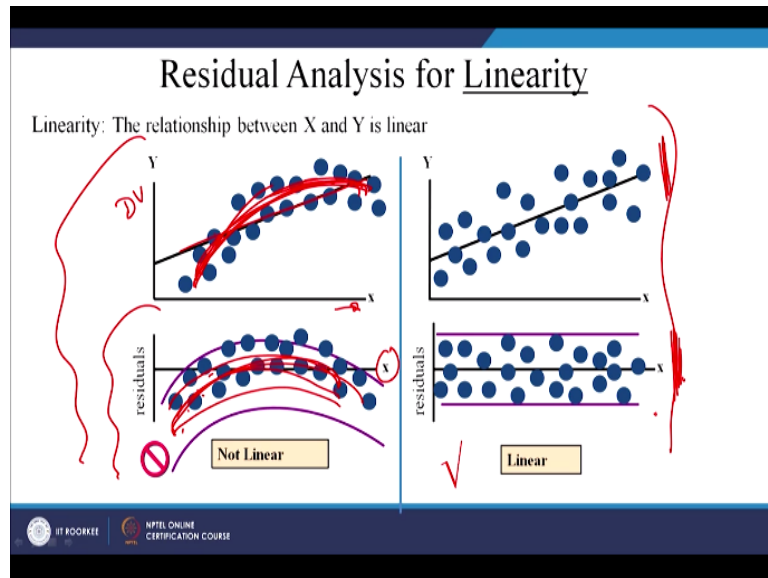
IT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE

Let us look at residual analysis in fact residual analysis are this error terms will help you in knowing whether you are fulfilling those assumptions or not. So observed and predicted value. So the differences between them, it can be positive or it can be negative. So check the

assumption of regression by examining the residuals or errors. So these are 4 assumptions which we have talked about right.

And you can get little idea about linearity of X and Y just by having plot of X and residual values right. So we can graphically see, it is not only linearity but independence of error term normality and equal variance okay.

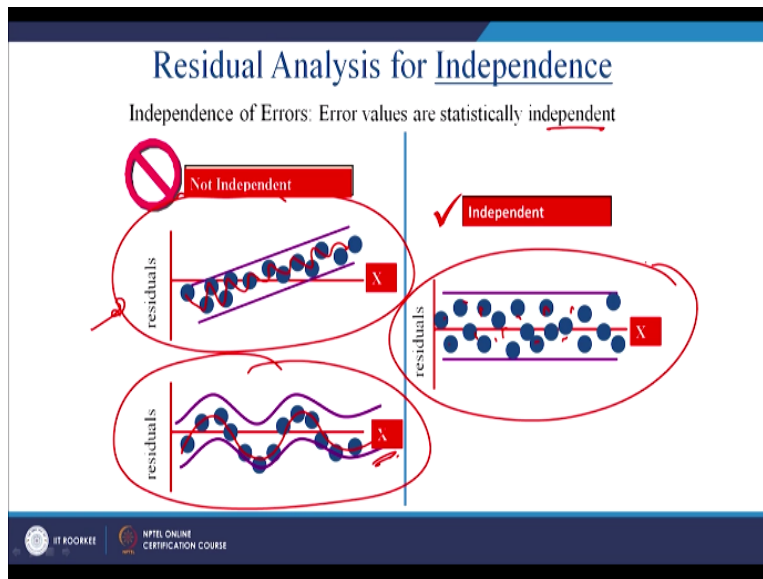
(Refer Slide Time: 04:50)



So this is how the relationship looks like, so you have got independent variable X, this is dependent variable Y, this is your if you see carefully this is not a linear relationship right, it is a curvilinear relationship and if you look at this it is relationship between X and residuals right. So you plot X on x-axis and residuals on y-axis. So if the residuals are like this, there is a pattern right or it is a nonlinear one right.

So it should not be like this, is not it? You should have linear data as well as the relationship between X and Y should be linear and the residual should be linear right with respect to X. So this is linearity assumption.

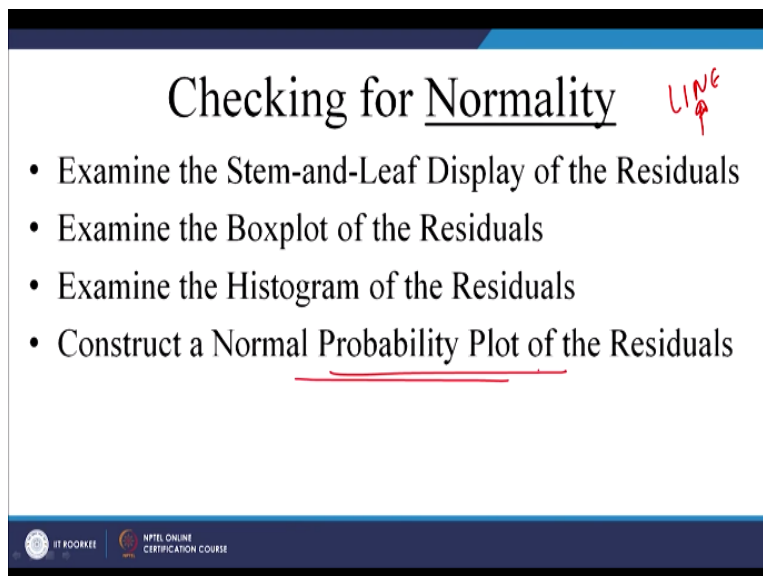
(Refer Slide Time: 05:51)



Independence of error terms so the error term should be independent of each other right, so let us look at this. So these are let us have one example, residuals versus X, again residual versus independent variable X. So if you look at carefully then there is a pattern over here, it not it. Up and down, up and down right similarly here, just see there is a pattern. So there should not be any pattern as far as residuals and X plot is concerned right.

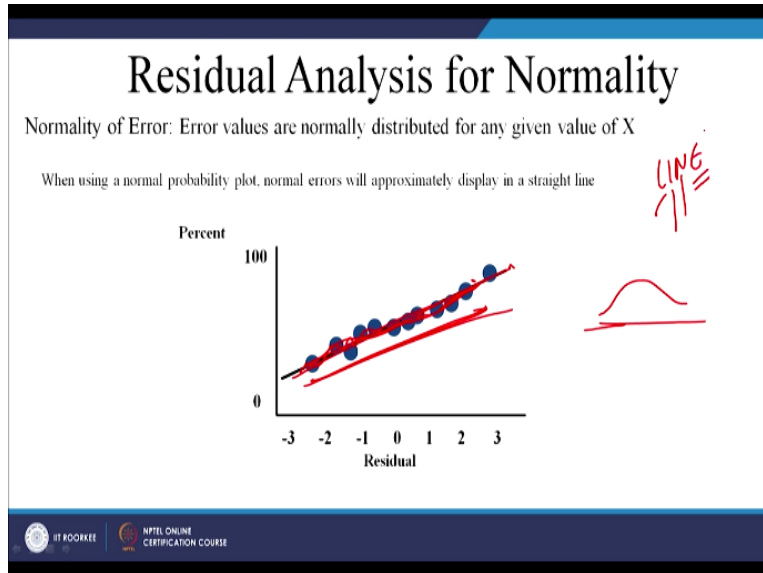
If there is a pattern, it means they are dependent on each other, so error terms should be independent of each other just like this. So here you are getting different values right. There is no pattern as such, so this is independence of error terms, second assumption.

(Refer Slide Time: 07:02)



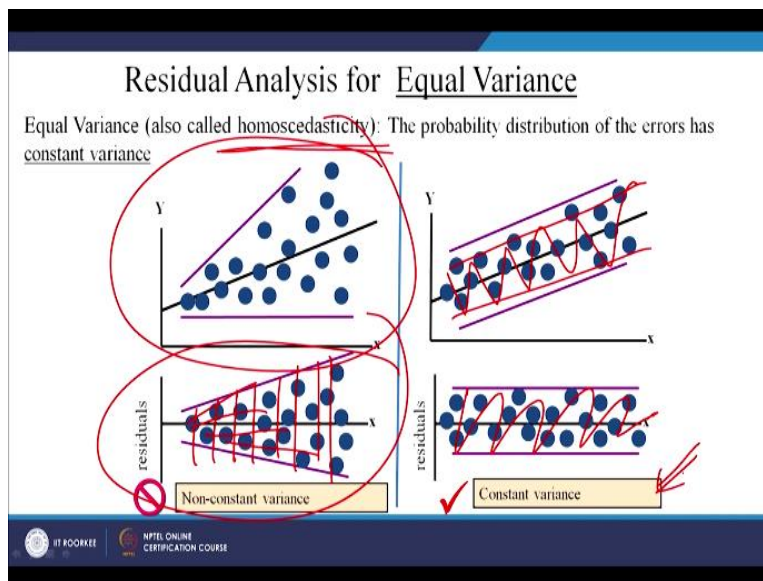
The third one is as I said LINE, so the third one is normality of data. So normality of data can be seen by multiple ways which we have seen in very few initial classes. So one of them is you know you can construct a normal probability plot.

(Refer Slide Time: 07:30)



So this is how you plot a normal probability plot. So if all these points are on the straight line, this is nothing but normal data right. In fact, this is nothing but your normal distribution right. So it is a bell-shaped curve but since you have taken log normal so that is why this is the shape of distribution. So we have seen 3 assumptions, LIN right linearity, independence of error terms, normality and the fourth is error term should be equally distributed right.

(Refer Slide Time: 08:06)



So equal variance, it is also known as homoscedasticity. So you should have homoscedasticity not the heteroscedasticity. So these are your error terms okay. So constant

variance is there and just see even in this plot as well as in this plot you have got this error terms are whereas these terms are spreading right. So this should not be there. It should be just equal, so this variance is constant here right, is not it? Similarly, here so the variance should be constant as far as error terms are concerned.

(Refer Slide Time: 09:02)

House Price in \$1000s (Y)	Square Feet (X)	Residuals = (observed - predicted) (Y) - \hat{Y}
245	1400	245 - 251.97 = -6.97
312	1600	38.12
279	1700	-5.86
308	1875	3.925
199	1100	-19.98
219	1550	-49.39
405	2350	48.77
324	2450	-43.21
319	1425	64.335
255	1700	-29.86

$$Y = a + bx = 98.25 + 0.1098 * x$$

$$= 98.25 + 0.1098 * 1400 = 251.97$$

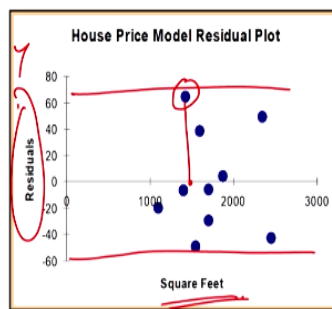
So let us find out those LINE assumptions and we will see how to get the residuals. So this is the question which we have worked out. So house price and the size of the plot and this was our line of best fit or prediction line right. So a is = 98.25 and b is 0.1098, so let us say what is the value of \hat{y} or predicted value when X is =1400 right. So in this equation just put all X values. So when you put X is =100, this \hat{Y} really becomes 259.

So $y - \hat{y}$ is residuals is not it? So residual over here is -6.97, this is plus, this is minus, this is plus, minus, minus, plus, minus, plus. So if you look at carefully there is a pattern, you have got plus, minus, plus, minus is not it? So these are residuals right and we will draw these residuals against X values okay.

(Refer Slide Time: 10:25)

Simple Linear Regression Example: Excel Residual Output

RESIDUAL OUTPUT		
	Square Feet (X)	Residuals
1	1400	-6.923162
2	1600	38.12329
3	1700	-5.853484
4	1875	3.937162
5	1100	-19.99284
6	1550	-49.38832
7	2350	48.79749
8	2450	-43.17929
9	1425	64.33264
10	1700	-29.85348



Does not appear to violate any regression assumptions

So this is how we have drawn these residuals which are here on y-axis and on x-axis you have got square feet right. So for X is =1400 you have got -6.92. So in fact you can see over here let us say the highest one is 64 right, 64.33 this is the one, 64.33 and for a given value of X is =1425 right, so this is 1425. So you can draw all these 10 points. Now if you look at all these 10 points, they are quite randomly distributed on plot.

We are not saying that even you can say once again here, so minus, plus, minus, plus. There is no such pattern right. If there is one negative value, the other is positive then negative, positive, negative to negative, positive, negative, positive and so on. So there is no such pattern over here. So will say that in fact you can see that there is equal variance so linearity is there, independence of error term is there, normality is there and error terms are having equal variance right. So we fulfill all these 4 assumptions right.

(Refer Slide Time: 11:57)

Inferences About the Slope

- The standard error of the regression slope coefficient (b_1) is estimated by

$$S_{b_1} = \frac{S_{YX}}{\sqrt{SSX}} = \frac{S_{YX}}{\sqrt{\sum (X_i - \bar{X})^2}}$$

where:

S_{b_1} = Estimate of the standard error of the slope

$S_{YX} = \sqrt{\frac{SSE}{n-2}}$ = Standard error of the estimate

SSE/n-2

IT KOOHKEE NPTEL ONLINE CERTIFICATION COURSE

Now let us look at inference about slope. So we have seen standard error, standard error is nothing but it was sum of square error divided by n-2 under root right. So standard error, so will have the standard error of regress slope coefficient. So we have got this slope of independent variable b1, so it will have some standard error as well right. So standard error can be calculated like this, so this is estimate of standard error of the slope.

So you know the standard error very well, you just divide it by this value, you will get standard error of slope right.

(Refer Slide Time: 12:54)

Inferences About the Slope: t Test

- t test for a population slope
 - Is there a linear relationship between X and Y?
- Null and alternative hypotheses
 - $H_0: \beta_1 = 0$ (no linear relationship)
 - $H_1: \beta_1 \neq 0$ (linear relationship does exist)
- Test statistic

$$t_{STAT} = \frac{b_1 - \beta_1}{S_{b_1}}$$

where:

b_1 = regression slope coefficient ✓

β_1 = hypothesized slope ✓

S_{b_1} = standard error of the slope ✓

d.f. = n - 2

y = mx + c

b1

(sketch of a bell curve)

IT KOOHKEE NPTEL ONLINE CERTIFICATION COURSE

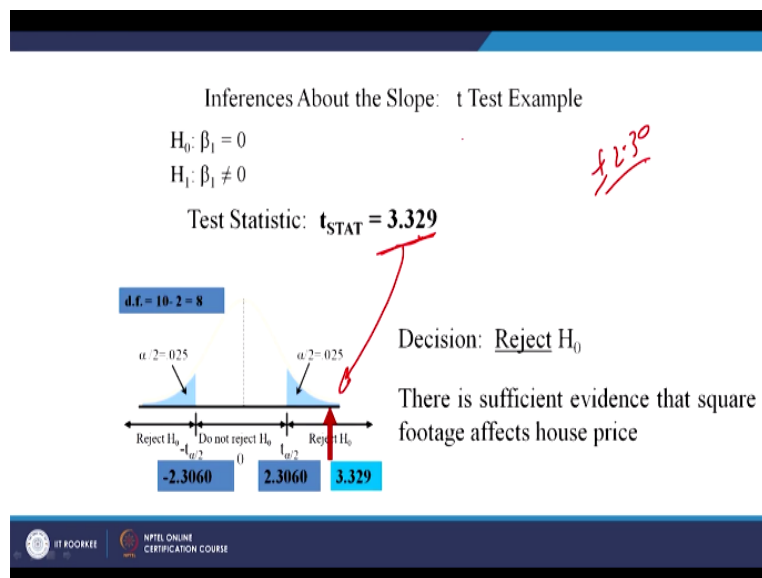
So inference about slope, now we want to know, you know we have seen $y = mx + c$, is not it? This is constant, this is slope m is slope and this is your independent variable right. Now in previous example, we have seen that the square feet, the size of the plot has effect on price

of the house, so we want to see whether this slope this m value or let us say b1 value or in case of multiple regression you can have b1, b2, b3 and so on.

So all these values have got significant effect on y or not, so is this effect is statistically significant or it is just by chance, so we have to check this relationship. So we say that the null hypothesis is or in other words we want to know is there any relationship between linear relationship between X and Y. So beta is 0, will say that there is no relationship. So if you reject null hypothesis then will say there is relationship right.

So initially we are saying there is no relationship right, so you can have just test statistics, so t statistics, you have got regression coefficient, hypothesized slope and standard error of slope at appropriate degrees of freedom. So calculate t-statistics value and since this is a case of two tailed test you will have two rejection regions, is not it? This is rejection region and this is second rejection region right. So this is first rejection region.

(Refer Slide Time: 14:49)



Now let us look at the output of the previous question and if you look at the t-statistics which we have calculated is 3.329 and the t-table value is ± 2.3 so ± 2.3 right. So calculated t-value is somewhere here right, so will reject null hypothesis. So null hypothesis is getting rejected right. So when we say null hypothesis is rejected, it means what it means there is linear relationship between X and Y right.

(Refer Slide Time: 15:29)

Inferences About the Slope: t Test Example

House Price in \$1000s (y)	Square Feet (x)
245	1400
312	1600
279	1700
308	1875
199	1100
219	1550
405	2350
324	2450
319	1425
255	1700

Estimated Regression Equation:

$$Y = a+bx = 98.25 + 0.1098 * \text{Sq feet}$$

The slope of this model is 0.1098

Is there a relationship between the square footage of the house and its sales price?

Now let us look at the regression line or estimated regression line for this question. So slope is this, is not it? So we have checked is there any relationship between X and Y or not. So though we have rejected the null hypothesis and we have said that there is relationship.

(Refer Slide Time: 15:51)

Inferences About the Slope: t Test Example

$H_0: \beta_1 = 0$
 $H_1: \beta_1 \neq 0$

From Excel output:

	Coefficients	Standard Error	t Stat	P-value
Intercept	98.24833	58.03348	1.69296	0.12892
Square Feet	0.10977	0.03297	3.32938	0.01039

From Minitab output:

Predictor	Coef	SE Coef	T	P
Constant	98.25	58.03	1.69	0.129
Square Feet	0.10977	0.03297	3.33	0.010

$$t_{STAT} = \frac{b_1 - \beta_1}{S_{b_1}} = \frac{0.10977 - 0}{0.03297} = 3.32938$$

The same information can be seen over here in output table either you are getting output from excel or from Minitab right. So just look at the P value, so P value is <0.05 so will reject the null hypothesis right. So we will say that there is relationship. In fact, the calculated t-value is available in table as well. So this is there 3.3293. Now standard error is there, of course all these values are there in this table.

And in fact output from Minitab also can be seen. Now P value again here is <0.05. So will say that will reject null hypothesis and there is relationship between X and Y. So far we have

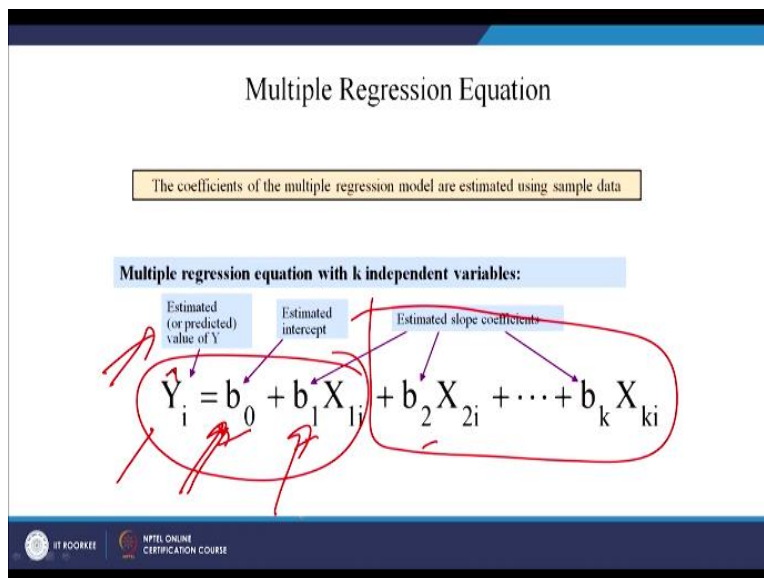
seen linear regression is called simple linear regression or bivariate regression because there were only two variables, one was dependent and the other one was independent variable and both of them were metric in nature. So we have seen simple linear regression right.

(Refer Slide Time: 17:10)



You may have a situation where the independent variables are more than 1, so in that case you need to apply something called multiple regression. So let us look at what is multiple regression.

(Refer Slide Time: 17:25)

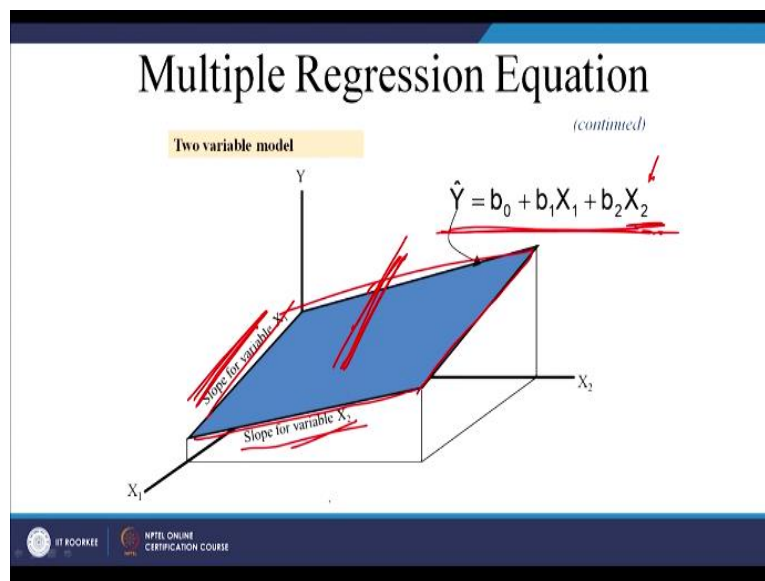


So it is very similar to simple linear regression, only the difference is you have got more than one independent variables right. So here you have got k independent variables, again you have got error term, so you have got plus and minus both in error term right. So this is so all

these are your independent variables and these are different slopes right, beta 1, beta 2, beta k and this is your intercept right.

What is intercept? Intercept is that value when all independent variables are 0 right. So this is how you can have the line of best fit or the estimation line. So y is $=b_0, b_1, b_2$ and so on right which is similar to simple linear regression. Only the thing is you have got these extra terms over here right. Earlier we had only these three right, dependent variable, intercept and slope right.

(Refer Slide Time: 18:36)




Now this is how the multiple regression or the plane of estimation, I will not call it line of estimation right. So you are getting a plane over here right, it is a two-dimensional plane right, is not it? Because there are two independent variables, so there are two slopes right. Slope for X_2 and this is slope for X_1 . So when there was one independent variable, it was a line, two independent variables a plane.



Suppose if there are 3 independent variables, what would be the shape of this graph or this plot? It will not remain a plane right, it would be a three-dimensional figure is not it? It could be cuboid rather than a plane if there are three independent variables.

(Refer Slide Time: 19:31)

Example: 2 Independent Variables

- A distributor of frozen dessert pies wants to evaluate factors thought to influence demand
 - Dependent variable: Pie sales (units per week)
 - Independent variables: Price (in \$) ✓
Advertising (\$100's) ✓
- Data are collected for 15 weeks



So let us look at this question wherein there are two independent variables. A distributor of frozen dessert pies wants to evaluate factors thought to influence demand. So the first factor is price, the second one is advertising. In other words, we want to know how these two independent variables namely price and advertising affecting pie sales right. So data were collected for last 15 weeks.

(Refer Slide Time: 20:13)

Pie Sales Example

Week	Pie Sales	Price (\$)	Advertising (\$100s)
1	350	5.50	3.3
2	460	7.50	3.3
3	350	8.00	3.0
4	430	8.00	4.5
5	350	6.80	3.0
6	380	7.50	4.0
7	430	4.50	3.0
8	470	6.40	3.7
9	450	7.00	3.5
10	490	5.00	4.0
11	340	7.20	3.5
12	360	7.90	3.2
13	440	5.90	4.0
14	450	5.00	3.5
15	360	7.00	2.7


Multiple regression equation:



$$\widehat{\text{Sales}} = b_0 + b_1 (\text{Price}) + b_2 (\text{Advertising})$$

$$\sum Y = b_0n + b_1 \sum x_1 + b_2 \sum x_2$$

$$\sum x_1 Y = b_0 \sum x_1 + b_1 \sum x_1^2 + b_2 \sum x_1 x_2$$

$$\sum x_2 Y = b_0 \sum x_2 + b_1 \sum x_1 x_2 + b_2 \sum x_2^2$$



And regression was done so you can in fact if you look at this particular question, you can solve this question manually using these 3 equations because this is a question where in this equation there are 3 unknown variables right b0, b1 and b2. So there are 3 unknown and 3 equations. So these are nothing but simultaneous equations you can solve for b0, b1 and b2.

(Refer Slide Time: 20:48)

Excel Multiple Regression Output

Regression Statistics					
Multiple R	0.72212				
R Square	0.52148				
Adjusted R Square	0.44172				
Standard Error	47.46341				
Observations	15				

Sales = 306.526 - 24.975(Price) + 74.131(Advertising)

ANOVA					
	df	SS	MS	F	Significance F
Regression	2	29460.027	14730.013	6.53861	0.01201
Residual	12	27033.306	2252.776		
Total	14	56493.333			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%
Intercept	306.52619	114.25389	2.68285	0.01993	57.58835	555.46404
Price	-24.97509	10.83213	-2.30565	0.03979	-48.57626	-1.37392
Advertising	74.13096	25.96732	2.85478	0.01449	17.55303	130.70888

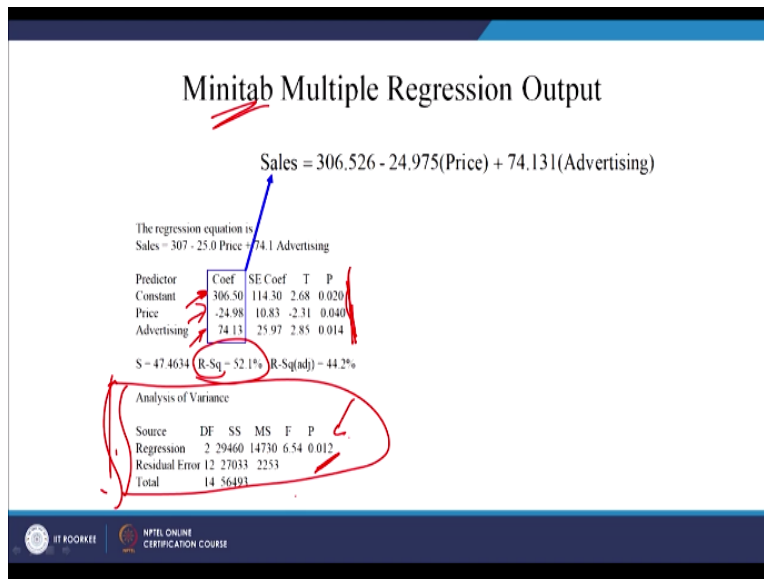
52.14%

However, you can solve this question using excel or Minitab as well right. So this is the output from excel. So if you look at the first part of this table where you get the descriptive statistics, so r square is 52.14% right. We will talk about adjusted r square little later. Standard error 47.46. How did we get standard error? Standard error is $Se = \sqrt{\text{sum of error term}/n-2}$ right and how did you get this r square?

R square is the regression sum of square right divided by total sum of square right SST, so this divided by this value you will get r square. So your estimation line is this, sales is = this is your intercept -24.97, this is coefficient of first independent variable, this is coefficient of second independent variable, so 74.13 and if you look at P values then all these 3 are significant, so will say that price and advertising as well as intercept especially price and advertising are significantly affecting dependent variable.

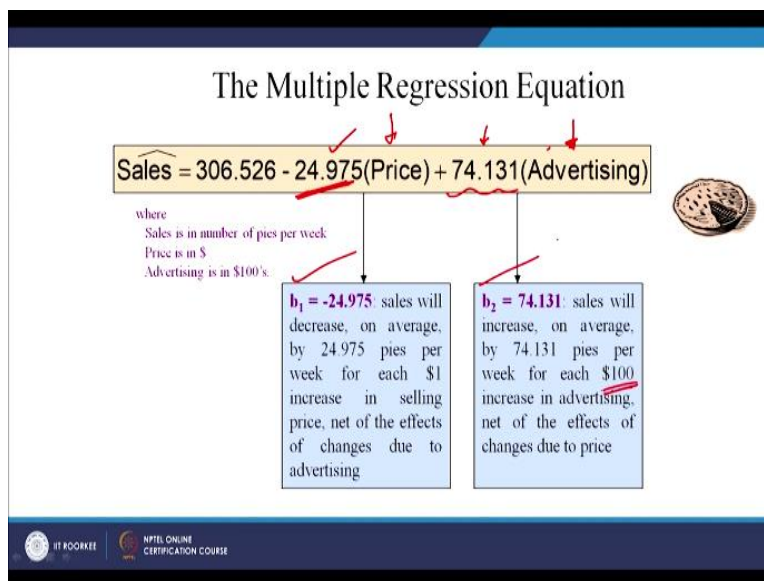
Look at this r square again, it is just 52.14% right. so will say that 52.14% variance independent variable is explained by these two independent variables. Look at ANOVA table, this part the middle part of the table ANOVA table, so this is nothing but the true value of F right so this is again < 0.05 so will say that the model is a significant model right.

(Refer Slide Time: 22:58)



So this is output from Minitab, will get the same output right, so constant or the intercept, price and advertising coefficients right. P value is again all are <0.05 , r square 52.1 same as this r square, then ANOVA table of course you can look at the how overall model is significant right. So this P is again it is <0.05 right.

(Refer Slide Time: 23:34)



So you can solve this question using Minitab but before this let me explain the meaning of these terms right. So this is -24.97 sales will come down by this much unit if you increase price of pie or this product by one minute. If you increase let us say advertisement, expenditure by 100 dollars then this much sales will increase, so this is the meaning of b_1 and b_2 right.

(Refer Slide Time: 24:12)

Using The Equation to Make Predictions

Predict sales for a week in which the selling price is \$5.50 and advertising is \$350?

$$\begin{aligned} \widehat{\text{Sales}} &= 306.526 - 24.975(\text{Price}) + 74.131(\text{Advertising}) \\ &= 306.526 - 24.975(5.50) + 74.131(3.5) \\ &= 428.62 \end{aligned}$$

Predicted sales is 428.62 pies

Note that Advertising is in \$100's, so \$350 means that $X_2 = 3.5$

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE

So predicted sales let us say if I ask you a question, what would be the sales when price is this much and advertising expenditure is this much, so put X1 is equal to this and X2 is equal to this, so this is X1 and this is X2 right. So will say that if the price is equal to this and advertising expenditure is equal to this, the predicted sales is 428.62 pies okay. So now we will work out this example.

(Refer Slide Time: 24:48)

Week	Pie Sales	Price	Advertising
1	350	5.5	3.3
2	460	7.5	3.3
3	350	8	3
4	430	8	4.5
5	350	6.8	3
6	380	7.5	4
7	430	4.5	3
8	470	6.4	3.7
9	450	7	3.5
10	490	5	4
11	340	7.2	3.5
12	300	7.9	3.2
13	440	5.9	4
14	450	5	3.5
15	300	7	2.7

Multiple Regression dialog box options:

- Confidence and Prediction Interval Estimates
- Confidence level for internal estimates: 95 %

Check the "confidence and prediction interval estimates" box

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE


And in fact you can use excel so it is very simple, you just put input, these variables as you know predictor variables and this as dependent variables right.

(Refer Slide Time: 25:09)

Coefficient of Multiple Determination

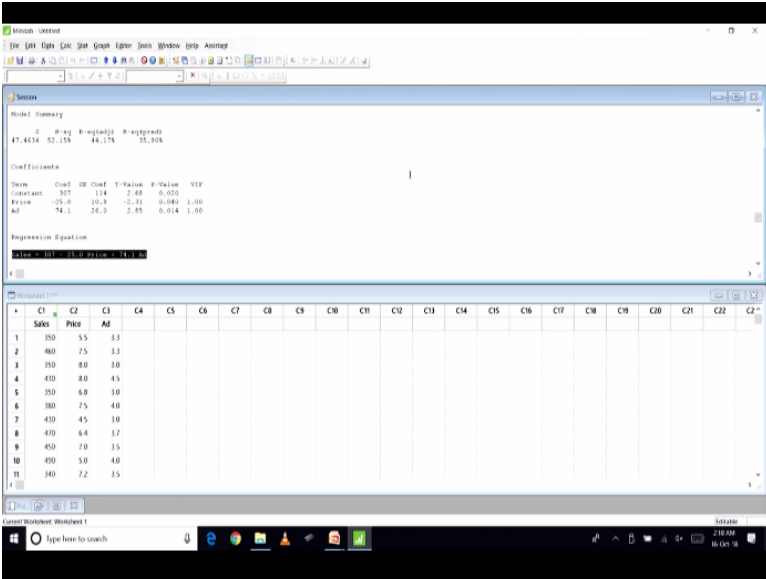
- Reports the proportion of total variation in Y explained by all X variables taken together

$$r^2 = \frac{SSR}{SST} = \frac{\text{regression sum of squares}}{\text{total sum of squares}}$$



So before moving onto next slide let me tell you how to solve this question using Minitab right. So this is the question we will solve using Minitab right. So pie sales is dependent variable this is X1, this is X2 right, two independent variables.

(Refer Slide Time: 25:26)



Row	Sales	Price	Ad
1	350	5.5	3.3
2	480	7.5	3.3
3	350	8.0	3.0
4	410	8.0	4.5
5	350	6.8	3.0
6	380	7.5	4.0
7	410	4.5	3.0
8	470	6.4	3.7
9	450	7.0	3.5
10	490	5.0	4.0
11	340	7.2	3.5

So let us look at this. So this is sales, this is price and this is advertising expenditure right, so 350, 480, 350, 430, 350, 380, 430, 470, 450 and 490. So there are some more data points 450 and 300. So there are total 15 data points. Similarly, for price now this 6 number, you have to enter it number 1, so just go up yeah so this is price 5.5, 7.5, 8, 6.8, 7.5, 4.5, 6.4, 7.0. So this is how you should enter all the prices 7.9, 5.9, 5 and 7.

Now you need to enter the advertising expenditure right and all these values are in terms of 100 right so it is 330, 100 no it is not 3.3, so the ninth one is 3.5. We are supposed to enter

data carefully otherwise there will be some difference in answer right. What I have shown you would be different and what we would be getting will be different. So we are supposed to enter data for this advertising expenditure.

Now go to stat, basic regression, yeah regression just regression, fit regression models right. So response is sales right sales then you have got continuous predictor right. So just click okay and let us look at output. So all these are P values right. So less than 0.05 will say these two are significant independent variables. Let us look at r value is 52.15%, then standard error is 47.46% and this is your line of best fit or estimation line right.

So $307 - 25 + 74.1$ which is exactly the same which you are getting over here right, this is 24.975 which is as good as 25. So this is how you can solve the question on multiple regression and if you have got let us say more than two independent variables, then we will have one more independent variable noted over here right with its positive or negative sign depending upon what is that independent variable.

So with this let me stop here. In next class, we will work out couple of examples on multiple regression and will see how to include a categorical independent variable in any problem. Thank you very much.