

Probability Methods in Civil Engineering
Prof. Dr. Rajib Maity
Department of Civil Engineering
Indian Institute of Technology, Kharagpur

Lecture No. # 40
Regression Analysis and Correlation (Contd.)

Hello and welcome to this lecture. Today's lecture is basically the continuation of the last lecture; we are discussing in this lecture about regression analysis and correlation. We have seen in the last class there are some of this basic correlation, basic regression analysis. We have seen its mainly that linear regression and those things after that we are in today's class we are going to cover that remaining part of this correlation analysis. First thing that we will go is that multivariate regression. In the last class, we have seen that in the simple linear regression, what you have seen? There is one input variable and other one is that independent variable, that is y versus x .

So that if the x is your input variable and y is your target variable or your independent variable then we have seen that how we can estimate the parameters for this regression model and we have done how, what are the original variance for the y and then after the regression what is the conditional variance and all.

Basically, we have also seen that through this regression what we are trying to do. How much, what is the extent of the variance of that target variable; that is, y is being reduced. So, herein this multiple regression what we do is that now this input variables are not 1, is more than 1. In that way we have to use the information of all those input variables and we have to develop a regression model for that target variable y . So, here the inputs are say x_1, x_2, x_3 and like that up to; obviously, x_m . So, there are, there could be some m variables, what should be the input?

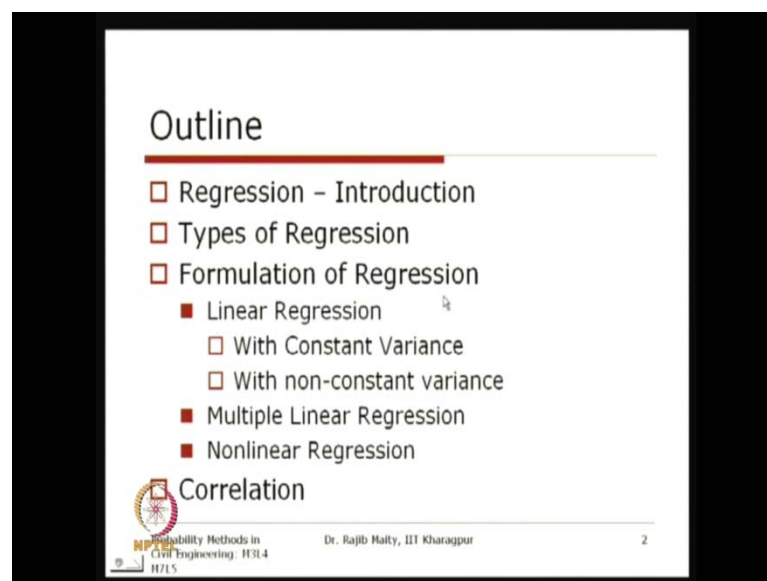
Now, if we just try to extend that analogy of this simple linear regression to this multiple linear regression, first we will take that then from this 1 dimensional to the 2 dimensional. As we are discussing in this last class that it is basically we are trying to take a straight line when it was a simple linear regression. But, in case now, if I just extend it instead of one input if it is two input then, basically it is a 3 dimensional space. That you

can imagine and through that 3 dimensional scatter plots of those points because, one point now will correspond to the three entries. One is that from x_1 other one is x_2 and the target is y .

So, basically one point consists of these three **a pair**. So, that the three entry, the three data point x_1 , x_2 and y . So, in the 3 dimensional scatter plot basically we are trying to fit a surface; now depending on if it is a linear regression then that surface will be a plane surface and now that surface should be the best fit to that through those points. And, through for that one following the same principle of this simple linear regression, we have to find out that it should be fitted in such a way that the sum of square error should be the minimum.

So, this is basic that transition from the simple linear regression to the multiple linear regression and based on this we will see what are the theories involved in this. So, we are continuing with that regression analysis and correlation.

(Refer Slide Time: 03:51)

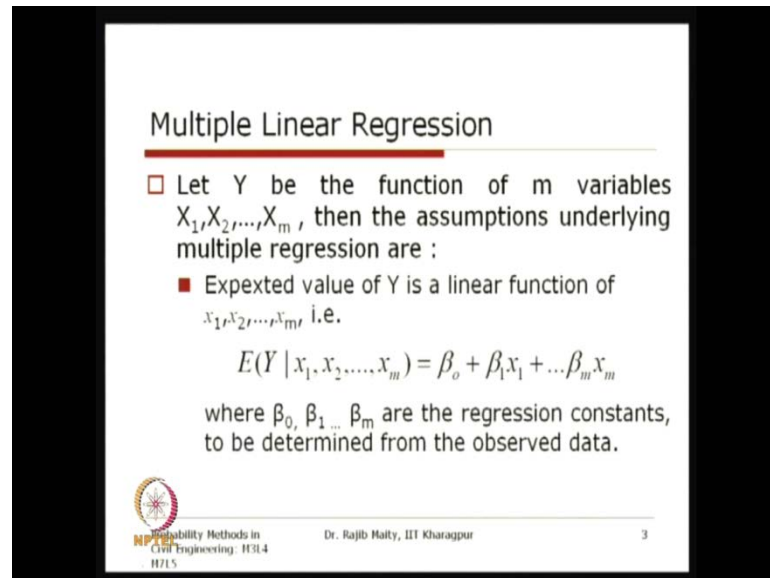


As we have seen in this last class, that in this different types of regression that we have discussed and there in the last lecture we have covered the linear regression and here we will see here today's class we will start with this multiple linear regression.

After we complete this one, we will see what is this non-linear regression and this non-linear regression can also be for both the cases. It may be before the simple linear regression

and also for the multiple linear regression and then we will go through this correlation and you know that this correlation we have discussed earlier also. But, here in the context of this regression analysis, we will once again see this aspect of this correlation. Basically, through this measure we are trying to identify what is the, how perfect the model that we have selected. So, that we will discuss under this correlation.

(Refer Slide Time: 05:01)



Multiple Linear Regression

- Let Y be the function of m variables X_1, X_2, \dots, X_m , then the assumptions underlying multiple regression are :
 - Expected value of Y is a linear function of X_1, X_2, \dots, X_m , i.e.

$$E(Y | X_1, X_2, \dots, X_m) = \beta_0 + \beta_1 X_1 + \dots + \beta_m X_m$$
 where $\beta_0, \beta_1, \dots, \beta_m$ are the regression constants, to be determined from the observed data.

NPTEL Probability Methods in Civil Engineering: H3L4 IITLS

Dr. Rajib Maiti, IIT Kharagpur

3

Well, to start with for this multiple linear regression, as I was just telling that there will be one target variable, we generally call it that independent variable, sorry, that dependent variable this is our target variable y and there are more than one independent variable which are X_1, X_2, X_m .

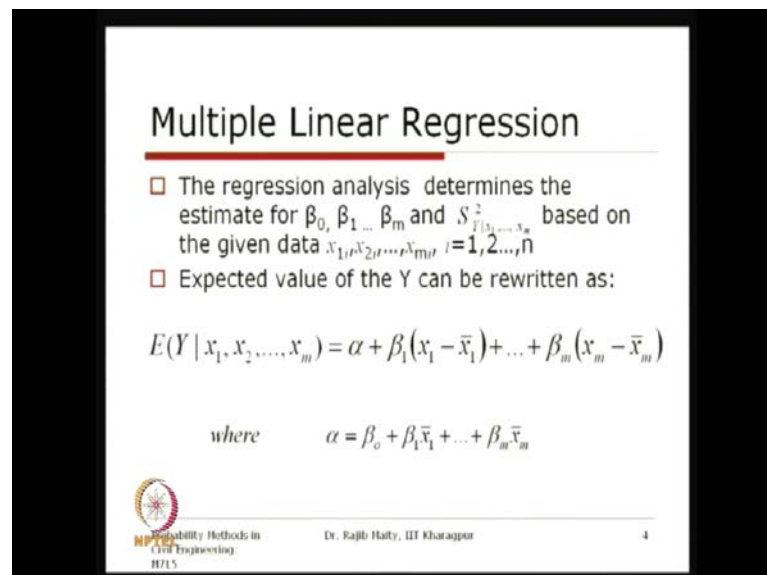
So, let Y be the function of m variables X_1, X_2 , up to X_m ; then the assumptions underlying the multiple regression are again following the same principle that we have discussed for the simple linear regression in the in the last lecture that we are trying to find out what is the expected value of this y given the input of this X_1, X_2, X_3 up to X_m .

So, this is the way that we are expressing that expectation of this target variable Y when the specific values of the input variables are given and this is expressed through this through this linear regression as we are now referring to this linear regression and that the linear regression is that β_0 plus $\beta_1 X_1$ plus $\beta_2 X_2$ up to plus $\beta_m X_m$.

If you recall, that in that last one we are having only one input and we are using this beta naught and beta 1X and there are only two parameters for this regression was there; beta naught and this one coefficient with that input variable. So, as you can see here there are m different inputs. So, this beta naught beta 1, beta 2, beta 3 up to beta m - these are our regression constants and this is to be determined based on the data that is available to us. You now can see that for this X1 if there are suppose that n numbers of that group data is available. Then, we are having that n numbers of X1 n numbers of X2 n numbers of X m.

So, through those points through those n points we have to fit one plane surface. Obviously, when we mention we are now referring to the context of these 2 input variable in a 3 dimensional case and obviously, that concept can be extended to the higher dimension. For example, here it will be them dimensional space that you can imagine. So, these constants are to be determined that is the basic underlying thing in this multiple linear regression as compared to the simple linear regression where the number of input was only one.

(Refer Slide Time: 07:53)



Multiple Linear Regression

- The regression analysis determines the estimate for $\beta_0, \beta_1, \dots, \beta_m$ and $S^2_{Y|X_1, \dots, X_m}$ based on the given data $X_{1i}, X_{2i}, \dots, X_{mi}, i = 1, 2, \dots, n$
- Expected value of the Y can be rewritten as:

$$E(Y | x_1, x_2, \dots, x_m) = \alpha + \beta_1(x_1 - \bar{x}_1) + \dots + \beta_m(x_m - \bar{x}_m)$$

where $\alpha = \beta_0 + \beta_1\bar{x}_1 + \dots + \beta_m\bar{x}_m$

Probability Methods in Civil Engineering IIT Kharagpur

Thus, this regression analysis determines the estimate for beta naught beta 1 up to beta m and that S y square that is the variance of this y given this X1, X2 up to X m based on the given data X1 i, X2 i up to X m i and i varies from 1 2 3 up to n. So, this i that the substitute that is used here that is basically represents the number of data that is available

to us and this 123 up to m represents this. So, this m represents that how many inputs that we are having now this S_y square given X_1, X_2, X_3 up to X_m is that the conditional variants of the target variable that is y.

So, this conditional variants should reduce with respect to this unconditional variants which is that S_y square and how much is this reduction that we can relate through there relate through that. So, more the reduction is better the model and at the end of this lecture as I was discussing, that at the end of this lecture, we will see that how that information is related to that correlation coefficient.

So, after some rearrangement of this equation that we have seen here that is expectation of y given this X_1, X_2, X_m equals to β_0 plus $\beta_1 X_1$ plus $\beta_2 X_2$ plus $\beta_m X_m$ this we can rearrange to get the another form of such of the same equation which is equals to the α plus β_1 into X_1 minus X_1 bar plus β_2 into X_2 minus X_2 bar like that up to β_m into X_m minus X_m bar.

So, this X_1 bar, X_2 bar or X_m bar is the mean of that particular input of that particular variable. So, as I was telling that X_i X_1 is that first input variable and this I can vary from one to m. So, we are having n data and this that mean of that input variable is represented this X_1 bar. Basically, this α what you can see now is basically one adjusted constant. Once again, including that β_0 and if you see from see that, α can be expressed like this; that α equals to β_0 plus $\beta_1 X_1$ plus $\beta_2 X_2$ bar up to $\beta_m X_m$ bar. So, this once we know this data basically this X_1 bar X_2 bar X_m bar are known and. So, once we get the estimate of this β_1 β_2 β_3 up to β_m and α with the help of that we can estimate that β .

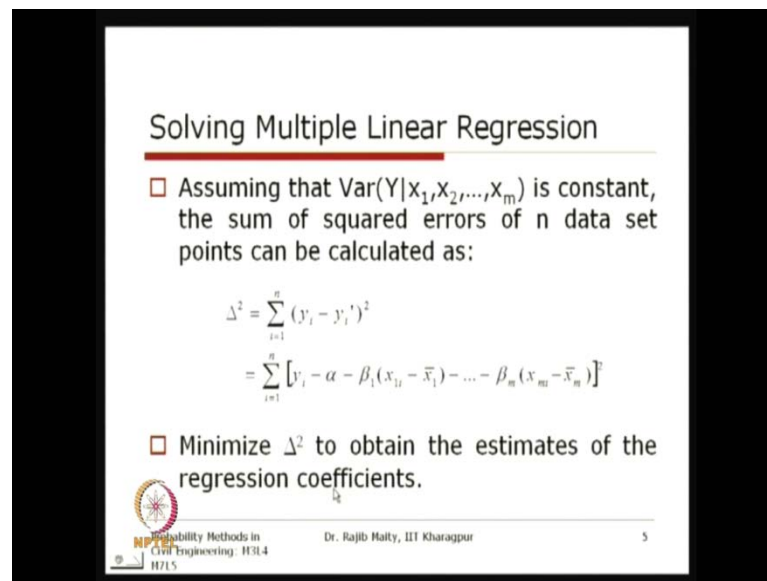
So, here in this expression what we can target is that we will first estimate these parameters α β_1 , β_2 , β_3 , β_m and with the help of this we will get what is the β_0 and we will get the final form of this regression equation like this.

So, now again we will use that same principle that we use there in the simple linear regression there basically as I was telling we are fitting a straight line and we are trying to minimize that error and that error means what is the error with respect to the modeled target variable and what is the observed. Here also, what we will do we will estimate these parameters in such a way. So that, that estimate of that of that variable Y and the observed Y their difference should be minimum. So, this difference now what

is actually observed and what is estimated from this model is basically your error and that error we should make it square and once we make that square and sum them up.

So, that is the sum of square error and with respect to that we will take the partial derivative with respect to all this constant that we are supposed to determine and that if we equate to 0 in the sense that we are minimizing that error sum of square error and we will get some simultaneous equation and we can solve this thing that we will see now.

(Refer Slide Time: 12:35)



Solving Multiple Linear Regression

- Assuming that $\text{Var}(Y|x_1, x_2, \dots, x_m)$ is constant, the sum of squared errors of n data set points can be calculated as:

$$\Delta^2 = \sum_{i=1}^n (y_i - y_i')^2$$

$$= \sum_{i=1}^n [y_i - \alpha - \beta_1(x_{i1} - \bar{x}_1) - \dots - \beta_m(x_{im} - \bar{x}_m)]^2$$

- Minimize Δ^2 to obtain the estimates of the regression coefficients.

NPTEL Probability Methods in Civil Engineering: H3L4 Dr. Rajib Maiti, IIT Kharagpur 5

So, now here again if we recall that again from this simple linear regression there are two cases. One is that the variants of Y with respect to that input variable whether that is constant or that can also vary. So, that in the case of this simple linear regression we have shown one case that where the variants can even vary with respect to the range of the input variable x .

Similarly, here also with respect to which zone we are talking about with respect to that the input variables. If the variants of this target variable that is Y is constant irrespective of this which zone that we are talking about, the combination of this input variable that is y that is X_1, X_2, X_m , then we can say that is either constant that variants is either constant or that can even vary. But, if it varies then that form that is how it is varying over this zone; so, that function should be known. So, now, if we assume that first case that is the conditional variants is constant then the sum of square error of the end data set points can be calculated as this.

So, this y_i is our actually observed target variable and y_i prime is that modeled variable that is we are getting from this regression equation. So, difference between them square them up and then sum up for all the individual observation that is n numbers of observations are there.

So, this quantity is giving that sum of square error now if we replace this one this y_i prime then it will come that this α minus β_1 into x_{1i} minus \bar{x}_1 minus β_2 into x_{2i} minus \bar{x}_2 minus like this up to β_m into x_{mi} minus \bar{x}_m and full quantity square. So, this will give you that sum of square error. Now, we have to minimize this Δ^2 to obtain the estimate for the regression coefficient.

(Refer Slide Time: 15:00)

Solving Multiple Linear Regression

i.e. $\frac{\partial \Delta^2}{\partial \alpha} = 2 \sum_{i=1}^n [y_i - \hat{\alpha} - \hat{\beta}_1(x_{1i} - \bar{x}_1) - \dots - \hat{\beta}_m(x_{mi} - \bar{x}_m)] = 0$

□ Similarly $\frac{\partial \Delta^2}{\partial \beta_1} = 0$ and $\frac{\partial \Delta^2}{\partial \beta_2} = 0$

□ From these set of equations we have:

$$\sum_{i=1}^n y_i - n\hat{\alpha} - \hat{\beta}_1 \sum_{i=1}^n (x_{1i} - \bar{x}_1) - \hat{\beta}_m \sum_{i=1}^n (x_{mi} - \bar{x}_m) = 0$$

$$\therefore \sum_{i=1}^n (x_{1i} - \bar{x}_1) = \dots = \sum_{i=1}^n (x_{mi} - \bar{x}_m) = 0$$

Thus $\hat{\alpha} = \frac{\sum y_i}{n} = \bar{y}$

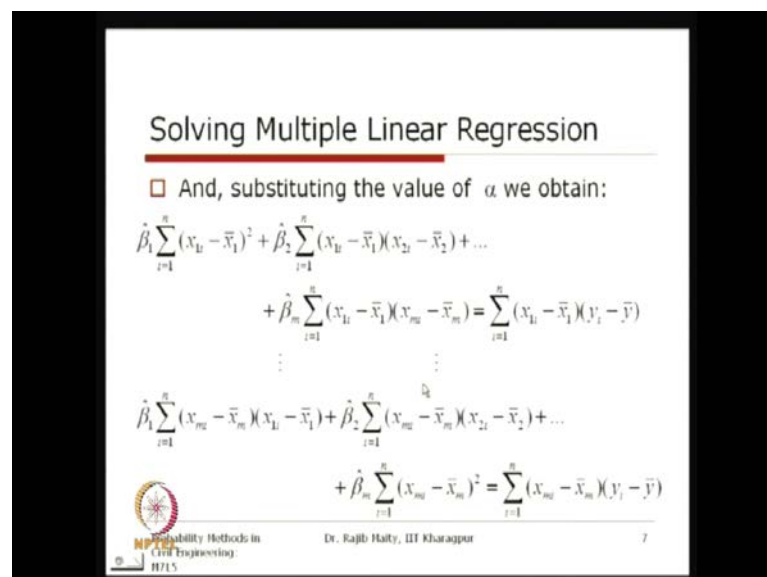
NPTV Probability Methods in Civil Engineering IIT Kharagpur 6

So, to minimize this one as we have seen in this simple linear regression also that with respect to all these constant we have to make it equal to 0 similarly that for this set Δ^2 this Δ^2 square α is equals to 0 similarly for this β_1 , β_2 , β_3 up to β_m like that all these partial derivatives should be equals to 0.

So, if you take this first one that Δ^2 square α and then we will we can take this one this partial derivative and that if we equate to 0 then the form comes like this where we can if we take this summation inside and we can see that this x_{1i} minus \bar{x}_1 bar. Basically, if we take this difference and sum them up; obviously, here power is 1. So, if we take that, sum them up, this will become 0.

So, like that for this x_2, x_3 up to x_m all these quantities will become 0 as it is written here. So, for all these quantities it will become 0. Now, if we just put this expression, these values here then, it will reduce to this form like this that α cap is equals to summation of y_i ; obviously, i from 1 to n divided by n . So, this is basically the mean of that observed target variable Y . So, this estimate of this α from this way we can see that it is the mean of Y now the remaining there are remaining n equations are there which are the partial derivative with respect to all n betas basically this β_1, β_2 up to β_m and there also we will get that α and this one this expression the estimate of this α if we just put it back.

(Refer Slide Time: 16:48)



Solving Multiple Linear Regression

□ And, substituting the value of α we obtain:

$$\hat{\beta}_1 \sum_{i=1}^n (x_{1i} - \bar{x}_1)^2 + \hat{\beta}_2 \sum_{i=1}^n (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2) + \dots$$

$$+ \hat{\beta}_m \sum_{i=1}^n (x_{1i} - \bar{x}_1)(x_{mi} - \bar{x}_m) = \sum_{i=1}^n (x_{1i} - \bar{x}_1)(y_i - \bar{y})$$

$$\vdots$$

$$\hat{\beta}_1 \sum_{i=1}^n (x_{mi} - \bar{x}_m)(x_{1i} - \bar{x}_1) + \hat{\beta}_2 \sum_{i=1}^n (x_{mi} - \bar{x}_m)(x_{2i} - \bar{x}_2) + \dots$$

$$+ \hat{\beta}_m \sum_{i=1}^n (x_{mi} - \bar{x}_m)^2 = \sum_{i=1}^n (x_{mi} - \bar{x}_m)(y_i - \bar{y})$$

NPTEL Probability Methods in Civil Engineering IIT Kharagpur 7

So, we can get a set of equations like this that this β_1 estimate of this β_1 this hat means here the estimate of this β_1 multiplied by $\sum_{i=1}^n (x_{1i} - \bar{x}_1)^2$ plus β_2 hat this $\sum_{i=1}^n (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2)$ like that all these plus and up to this β_m which is equals to that particular variable $\sum_{i=1}^n (x_{1i} - \bar{x}_1)(y_i - \bar{y})$ multiplied by $y_i - \bar{y}$.

So, this summation for these all in observation and so, this is the first equation. Like that, we can get all other equations and up to that m -th equation will be like this that β_1 cap into $\sum_{i=1}^n (x_{mi} - \bar{x}_m)(x_{1i} - \bar{x}_1)$ like that for this β_2 and β_m it is $\sum_{i=1}^n (x_{mi} - \bar{x}_m)(x_{mi} - \bar{x}_m)$ whole square which is equals to $\sum_{i=1}^n (x_{mi} - \bar{x}_m)(y_i - \bar{y})$ sorry $y_i - \bar{y}$ multiplication summation for 1 to n . One thing that you can see here

when we are taking the partial derivative with respect to β_1 , that is the first equation, this quantity is becoming square. Basically, this quantity is that x_{1i} minus \bar{x}_1 is basically getting multiplied for all the left hand side of this equation for all other variables.

So, these things x_{1i} minus \bar{x}_1 this one is there for all these entry and here also right side also you can see that this is the target this is related to the target variable Y and this one is that related to that for which constant we are taking the partial derivative here it is for β_1 . So, we can say that this quantity is this. Like that for this last one which is the m -th equation the partial derivative taken with respect to the β_m . So, you can see that what we have seen for the first expression is here that is x_{mi} minus \bar{x}_m square. So, this quantity is multiplied with this all other expression in the left hand side and also on the right hand side with which this function is related to the target variable which is also multiplied by this x_{mi} minus \bar{x}_m .

So, similarly you can see that for other variables also. So, if it is for this β_2 then the first quantity will be that x_{1i} minus \bar{x}_1 multiplied by x_{2i} minus \bar{x}_2 and the second quantity will be x_{2i} minus \bar{x}_2 whole square and like that right hand side will be x_{2i} minus \bar{x}_2 multiplied by y_i minus \bar{y} . So, in this way these are the simultaneous m equations are there and there are m unknowns.

(Refer Slide Time: 19:44)


Solving Multiple Linear Regression

□ Thus, we have 'm' linear simultaneous equations with 'm' unknowns, which can be solved for the values of the coefficients β_i and obtain the least squares regression equation

$$E(Y | x_1, x_2, \dots, x_m) = \hat{\alpha} + \hat{\beta}_1(x_1 - \bar{x}_1) + \dots + \hat{\beta}_m(x_m - \bar{x}_m)$$

$$= \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_m x_m$$

where $\hat{\beta}_0 = \hat{\alpha} - \hat{\beta}_1 \bar{x}_1 - \dots - \hat{\beta}_m \bar{x}_m$



Probability Methods in
Civil Engineering: H3R.4
H3L5

Dr. Rajib Maitty, IIT Kharagpur

8

Thus, we have the m linear simultaneous equation with m unknowns which can be solved for the values of the coefficient beta i. So, beta 1, beta 2, up to beta m and obtain the least squares regression equation. Once we get this estimate finally, the expression that linear regression expression that we are getting is that alpha cap plus beta 1 cap into x 1 minus x 1 bar plus up to this beta m hat means that the estimate of beta m into x m minus x m bar and now if we write in terms of this first equation that is the beta naught hat plus beta 1 hat x 1 plus beta 2 hat x 2 plus beta m hat x m.

This beta naught is now that beta the beta naught hat that is the estimate of this beta naught equals to alpha hat estimate of this alpha minus beta 1 hat x 1 bar minus beta 2 hat x 2 bar minus beta m hat x m bar. So, this is the final expression for this multiple linear regression that we will get.

(Refer Slide Time: 20:57)

Solving Multiple Linear Regression

- The conditional variance is calculated by :

$$S_{Y|X_1, \dots, X_m}^2 = \frac{\Delta^2}{n - m - 1} \quad (\text{unbiased estimate})$$

$$= \frac{\sum_{i=1}^n \left[y_i - \hat{\alpha} - \hat{\beta}_1(x_{i1} - \bar{x}_1) - \dots - \hat{\beta}_m(x_{im} - \bar{x}_m) \right]^2}{n - m - 1}$$
- And so, the corresponding standard deviation is obtained as :

$$S_{Y|X_1, \dots, X_m} = \frac{\Delta}{\sqrt{n - m - 1}}$$
- where n is the sample size and m is the number of dependent variables

NPTEL Probability Methods in Civil Engineering: H3R.4
Dr. Rajib Maiti, IIT Kharagpur

The conditional variance given that all those input variables will be that delta square this is the sum of square error divided by n minus m minus 1 which you can see that this delta square is the summation of this y i minus this alpha hat minus beta 1 cap. So, this is the sum of square error divided by n minus m minus 1.

So, now you recall that when there is only one input was there in case of the simple linear regression if you recall that equation for this conditional variance was delta square divided by n minus 2. So, there m was only 1. So, one input was there. So, if you put m equals to 1 here you will get the same expression that is delta square divided by n minus

2. So, if m increases, here m is more than 1. You have to put that value and this n is this number of data points that is available. So, this is one unbiased estimate why you have already discussed earlier that unbiased estimate. So, that number of degrees of freedom is lost; that should be reduced.

So, here you see that $n - m$. So, there are m estimates are there for this x_1 bar, x_2 bar up to x_m bar this is the estimate that we are there are those estimates are there. So, m degrees of freedom is lost here again one more parameter here the α hat is there that is that we are calculating.

So, it is that basically $m + 1$ degrees of freedom is lost and in case of simple linear regression, there are two constants were there that is two estimates. So, 2 degrees of freedom was lost there. So, it was Δ^2 divided by $n - 2$ in that simple linear regression. So, this is the unbiased estimate of the conditional variance of y .

Similarly, if that corresponding standard deviation is the positive square root you know. This is Δ by square root of $n - m - 1$ where n is the sample size and m is the number of dependent variables. So, now we will take up one problem here. So, you can see that there are many applications can happen in the civil engineering problem.

(Refer Slide Time: 23:26)

$m=2$

Observation No.	Y	X1	X2	$(Y - \bar{Y})^2$	$(X1 - \bar{X1})^2$	$(X2 - \bar{X2})^2$	$(Y - \bar{Y})(X1 - \bar{X1})$	$(Y - \bar{Y})(X2 - \bar{X2})$	$(X1 - \bar{X1})(X2 - \bar{X2})$	$(Y - \bar{Y})$
1	45.2	2210	38.20	394258	1.345	671.853	-4081	-6.96	49.26	16.472
2	30.1	704	35.10	771060	4.121	1782.541	1405	3.25	55.34	27.452
3	53.5	1349	39.30	54336	4.709	-505.827	-420	3.91	51.52	3.910
4	45.3	3528	36.10	3786527	1.081	-2004.277	-12454	6.59	45.97	0.447
5	50.2	2550	34.50	918030	6.917	-2545.577	-1452	3.95	49.74	0.209
6	57.5	955	38.30	393254	1.369	-733.707	-3637	6.79	53.19	18.549
7	53.2	1980	35.70	158124	2.045	-568.997	597	-2.15	51.04	4.653
8	59.3	650	33.10	868810	16.241	-3756.363	-7084	-30.63	56.35	8.673
9	48.2	645	39.80	878156	7.129	-2502.057	3280	-9.34	53.54	28.540
10	54.5	2250	41.20	110290	16.565	-1351.647	-4920	11.40	51.03	12.012
Sum	517.000	13421	372.300	8351847	81.301	6062.880	-24776	-13.200	120.318	
Mean	51.700	1342.1	37.230	835185						

$\alpha = 72.4$
 $\text{var}(Y) = 22.78$
 $\beta_1 = -0.422$
 $\beta_2 = -0.0032$
 $\beta_3 = 72.4$

$E(Y|X_1, X_2) = 72.4 - 0.422x_1 - 0.0032x_2$
 $S^2_{Y/X_1, X_2}$

Cond Var = 17.274
 Cond Std = 4.156
 Fsq = 0.242

If you just see here for this particular problem here, this could be any say for example, that we are talking about what should be the temperature with respect to its altitude and

latitude generally the altitude increases. So, we know that the temperature will decrease and if the latitude also increases the temperature generally decreases.

So, suppose that type of data if we just take here. So, this is the number of observation that 10 numbers of observations are given and this is the target variable of this Y. So, you can see that it is 45.25.1 like this. So, these are the ten observations that you can see here and these are the inputs of this see here we have taken the two inputs only x_1 and x_2 . So, here m is equal to your 2. I can write that m is equal to 2 here. So, this is your x_1 and this is your x_2 . So, this is up to this you can see that this is that data that we are getting.

Now, to get that estimate if you want to know that. So, these three columns that you can see is the input. So, first thing that you have to calculate is that this could be used in a general trade sheet; just to explain that how these things can happen. First what we can calculate, that this x_1 , x_1 minus \bar{x}_1 that mean that we calculated here. So, that square up for this one. So, basically this x_1 minus this mean is calculated here. So, this minus this that square will be give you this value and similarly, for this all such values we can calculate this one.

Similarly, this is for the x_2 minus \bar{x}_2 that square that you can calculate and then this is your x_1 minus \bar{x}_1 mean multiplied by x_2 minus \bar{x}_2 minus that \bar{x}_2 mean and their multiplication; it is basically is that column end. So, this entry minus this multiplied by this entry minus this mean. So, this is your mean row that you can see here and this is the summation and there are ten observations are there. So, this 2 to 10 minus this value multiplied by this 38.2 minus this 37.13 this is the mean if you multiply this we will get this one.

Similarly, this column is for that x_1 minus \bar{x}_1 mean multiplied by y , y minus \bar{y} mean. So, if you see this one here also basically, this minus this mean multiplied by this minus this mean we will get these values. So, here again this column if you see that then this y_2 minus \bar{y}_2 multiplied by y_i minus \bar{y} . So, this x_2 minus this \bar{x}_2 mean into that y minus \bar{y} that multiplication if you take we will get this value similarly for all these ten observations we have calculated and the last row that you can see is their summation.

So, up to this we can just calculate first directly based on whatever the data that we are having and then you know that estimate of this alpha is equal to your that y mean. So,

this is directly the 51.7 that we have seen for the mean for this Y now the variants of Y also you can calculate whatever the we have seen this the that Y that we can calculate the variants of that from this sample estimate that we discuss earlier you will get this 22.78.

Now, this expression if you refer to this equation here that is that; just here you have to put that m equals to 2. So, this value minus this square multiplied by this β_2 into this one this value will be equal to this one. So, all these quantities have calculated now here. Using this information, these quantities that just now what we have calculated we will just set. So, here there will be two simultaneous equations. So, these two equations are written here. So, this is your that first quantity multiplied by β_1 plus this quantity multiplied by β_2 is equals to the right hand side.

Similarly, this is the first equation; this is the second equation that we get and if we solve these two equations and there are two unknowns β_1 and β_2 . So, you will get that this is your β_1 and this is your β_2 β_1 is minus 0.0032 and β_2 is minus 0.422. So, with this estimate and that α also we know. So, that β_0 we can calculate that is β_0 as you have seen here that β_0 is equals to your α cap minus β_1 hat $\times \bar{x}_1$ minus β_2 hat $\times \bar{x}_2$ bar.

So, like if you use that one. So, β_1 minus this mean of this two that α minus β_1 estimate multiplied by the mean of x_1 minus β_2 estimate into that \bar{x}_2 then we will get what is your β_0 . So, once we get these three information that is your this is your β_0 this is β_1 and this is β_2 then we are getting this full expression like this that Y is equals to your 72.4 which is your β_0 minus 0.422 this is the estimate for this β_1 multiplied by x_1 minus 0.0032 $\times x_2$ is that expression.

Now, if we use this expression and use this what is this input then what we are getting this is the estimate of this Y basically you can say that is the expectation of Y given x_1 and x_2 equals to this. So, this quantity is basically is calculated in this column. So, which is the estimate for this Y? So, I can write that y' and these are the error. So, this Y minus y' , y' ; these are the error and that square. So, if I just square for all these 10 observations and sum them up this is basically is your that capital delta square that we are getting here.

So, once we get that capital delta square then we can calculate what is the conditional variants. This conditional variants means as we have seen that S_y is square given that x_1 comma x_2 . So, this if we calculate we will get this 17.274 and the conditional standard deviation is this positive square root of this one which is 4.156 and then square value that we get is that 1 minus.

So, how much is this reduction? So, you can see that these variants of this Y that is unconditional variants was 22.78 and the conditional variant is 17.274. So, that we can see here and this that r square is the. So, that 22, sorry 24.2 percent has is the reduction.

(Refer Slide Time: 31:36)

	Y	X1	X2	$\sum Y^2$	$\sum X1^2$	$\sum X2^2$	$\sum YX1$	$\sum YX2$	$\sum X1X2$
1	45.7	2710	38.20	2087.49	7356.10	1459.24	1517.89	-4001.33	-13200
2	50.1	2704	35.10	2510.01	7310.82	1232.01	1355.81	-4205.92	-12454
3	53.5	1340	38.30	2862.25	1795.60	1466.49	4799.05	-4096.92	-4205.92
4	45.3	2520	36.10	2051.81	6350.40	1303.21	1061.05	-3004.27	-12454
5	50.2	2650	34.50	2520.04	7022.50	1190.25	1261.00	-3545.57	-14852
6	51.5	950	28.20	2652.25	902.50	795.24	1308.75	-1221.07	-2627
7	53.2	1800	35.70	2830.24	3240.00	1274.49	2045.60	-588.87	-657
8	58.2	550	22.10	3387.24	302.50	488.41	1024.20	-304.04	-2052
9	48.2	845	38.80	2321.64	714.02	1505.44	7139.28	-2822.08	-3300
10	54.5	1250	41.20	2970.25	1562.50	1697.44	10565.00	-1251.84	-800
11	517.000	15671	377.300	267289.00	24546.41	14182.09	81301.00	-4001.330	-13200
12	517.000	15671	377.300	267289.00	24546.41	14182.09	81301.00	-4001.330	-13200
13	517.000	15671	377.300	267289.00	24546.41	14182.09	81301.00	-4001.330	-13200
14									
15									
16	41.70	8381947	10481	1740330	7028300	10983	184778	-34778	-13200
17	22.78	4001.230	10481	81301	7028300	10983	184778	-34778	-13200
18									
19									
20									
21									
22									
23									
24									
25									
26									
27									
28									
29									
30									
31									
32									
33									
34									
35									
36									
37									
38									
39									
40									
41									
42									
43									
44									
45									
46									
47									
48									
49									
50									
51									
52									
53									
54									
55									
56									
57									
58									
59									
60									
61									
62									
63									
64									
65									
66									
67									
68									
69									
70									
71									
72									
73									
74									
75									
76									
77									
78									
79									
80									
81									
82									
83									
84									
85									
86									
87									
88									
89									
90									
91									
92									
93									
94									
95									
96									
97									
98									
99									
100									

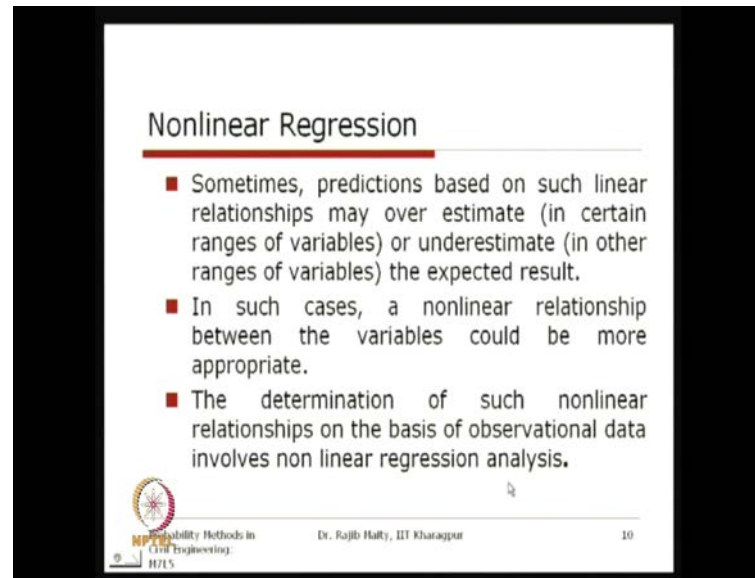
If you want to see this one in this larger font, you can refer to this excel file which you can see here more clearly whatever the calculation that we have done.

So, with this we can say here once again that now, when should I say that this model that what we have whatever we have we got that is strong enough or not? That we will see that in terms of this correlation coefficient that we are going to discuss now and again this will be, this is basically a part of this hypothesis testing.

Once we estimate one parameter and that parameter whether that is significant or not that can be tested through this hypothesis testing that we have covered earlier. So, through that


hypothesis testing we can see that how significant that estimate is and the hypothesis testing was covered in the earlier lecture.

(Refer Slide Time: 32:41)



Nonlinear Regression

- Sometimes, predictions based on such linear relationships may over estimate (in certain ranges of variables) or underestimate (in other ranges of variables) the expected result.
- In such cases, a nonlinear relationship between the variables could be more appropriate.
- The determination of such nonlinear relationships on the basis of observational data involves non linear regression analysis.

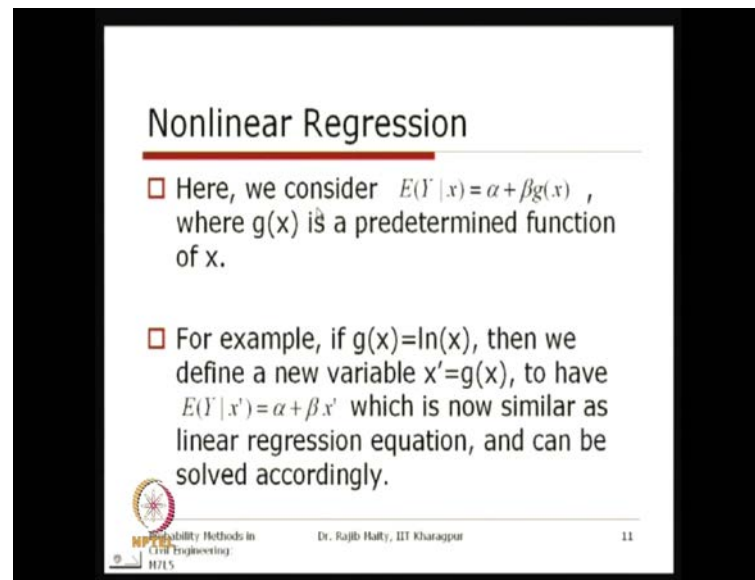
 Probability Methods in Civil Engineering
IIT Kharagpur

Dr. Rajib Halty, IIT Kharagpur

10

So, here we will now we will go to that non-linear regression and this non-linear regression is essential when we see that when that linear in case of this simple linear regression. We generally get one straight line and for this if there are two inputs we get a straight surface. But, if that straight line or the surface which is linear in nature may not express the variability fully, what we get sometimes is that the predictions based on such linear relationship may over estimate in certain ranges of this variable or underestimate in other ranges of this variable of this expected result. In such cases, a non-linear relationship between the variables could be more appropriate the determination of such non-linear relationship on the basis of the observational data involves a non-linear regression analysis.

(Refer Slide Time: 33:38)



Nonlinear Regression

- Here, we consider $E(Y|x) = \alpha + \beta g(x)$, where $g(x)$ is a predetermined function of x .
- For example, if $g(x) = \ln(x)$, then we define a new variable $x' = g(x)$, to have $E(Y|x') = \alpha + \beta x'$ which is now similar as linear regression equation, and can be solved accordingly.

Probability Methods in Civil Engineering
PMCE
2015

Dr. Rajib Hatty, IIT Kharagpur

11

So, basically in this non-linear regression what we do is that the expectation of y on condition x is equals to a function like this $\alpha + \beta g(x)$. So, this $g(x)$ is a predetermined function of this x . What we can do is that whatever the input x is there we will transfer; we will get a another new variable through this $g(x)$ and that that variable I can use with respect to this y and follow again that either the simple linear regression or multiple linear regression because, once we have converted it then you can see this form of this equation is a linear regression form.

Say for example, that if that $g(x)$ is equals to that log natural of this x then we can define a new variable x' of $g(x)$ to have that expectation of y given x' . Now, it is the converted variable is equals to $\alpha + \beta x'$ which is now similar as the linear regression equation and can be solved accordingly.

(Refer Slide Time: 34:43)

Example

Q. The average all-day parking cost in various cities of India is expressed in terms of the logarithm of the urban population, that is modelled with the following nonlinear regression equation:

$$E(Y|x) = \alpha + \beta \ln x$$

with a constant $\text{Var}(Y|x)$, where

Y = average cost in Indian Rupees for all-day parking cost (in Hundreds)

x = urban population (in thousands)

Estimate $\alpha, \beta, \text{Var}(Y|x)$ on the basis of the observed data.

Probability Methods in Civil Engineering
Dr. Rajib Hatty, IIT Kharagpur
12

We will take up one example. The average all day parking cost in various cities of India is expressed in terms of the logarithm of the urban population that is modeled with the following non-linear regression equation. The expectation of Y given x is equal to your alpha plus beta ln x with a constant variance Y given x where Y is the average cost in Indian rupees for all day parking cost in hundreds and x is the urban population in thousands.

This relationship is sometimes what happens if we just plot the data? That is, say, for example, here the Y and x if we plot it through a scattered plot that time that is nature can be visible whether a linear or a non-linear expression should be more appropriate or not. So, these are some initial guess. So, based on that, if we see that this log transform of that x might be the better estimate for this case, that is why the proposed equation is that alpha plus beta of log natural of that of x. So, we have to estimate that alpha and beta.

As we have told that we have to first transform that values of this urban population through this log natural function and then we will get a new set of this new variable, new observation in place of this x and then we can follow whatever we have seen in the linear regression equation.

(Refer Slide Time: 36:27)

Example...Contd.

City	x_i	y_i
1	300	0.51
2	280	0.47
3	330	0.57
4	450	0.59
5	370	0.65
6	540	0.83
7	450	0.87
8	1990	0.96
9	3360	1.50
10	3560	1.13

Dr. Rajib Halty, IIT Kharagpur

So, here there are 10 cities; this expression that is the x_i that in thousands that is what is the population that is shown here and this is the y_i in hundreds; what is the parking cost for all day parking?

(Refer Slide Time: 36:49)

Example...Contd.

City	x_i	y_i	$x_i' = \ln x_i$	$x_i' y_i$	$(x_i')^2$	y_i^2	$y_i' = a + \beta x_i'$	$(y_i - y_i')^2$
1	300	0.51	5.704	2.909	32.533	0.260	0.563	0.003
2	280	0.47	5.635	2.648	31.751	0.221	0.543	0.005
3	330	0.57	5.799	3.305	33.629	0.325	0.591	0.000
4	450	0.59	6.109	3.604	37.323	0.348	0.681	0.008
5	370	0.65	5.914	3.844	34.970	0.423	0.624	0.001
6	540	0.83	6.292	5.222	39.584	0.689	0.734	0.009
7	450	0.87	6.109	5.315	37.323	0.757	0.681	0.036
8	1990	0.96	7.596	7.292	57.698	0.922	1.114	0.024
9	3360	1.50	8.120	12.180	65.929	2.250	1.266	0.055
10	3560	1.13	8.178	9.241	66.872	1.277	1.283	0.023
			8.080	65.454	55.560	7.471		0.164

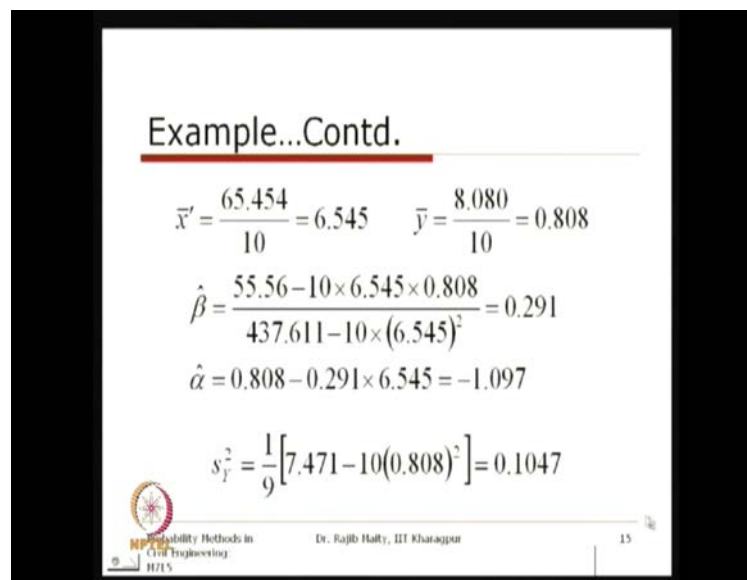
Dr. Rajib Halty, IIT Kharagpur

So, with these two data what we can first do is that, this x_i that is the input that is converted through this log natural and we get this expression. So, using this one as the input, we have to model this y_i . So, all those quantities that we require for this least

square estimate that is that should be estimated only in place of this x_i we should use this x_i' .

So, x_i' multiplied by y_i then x_i' square then y_i square this we can calculate. So, up to this column whatever the data that is available, we can calculate and we can take their individual summation also. Now, with the help of this information this summation - the summation of x_i' summation of x_i' square summation of y_i summation of y_i square; we can get that.

(Refer Slide Time: 37:42)



Example...Contd.

$$\bar{x}' = \frac{65.454}{10} = 6.545 \quad \bar{y} = \frac{8.080}{10} = 0.808$$

$$\hat{\beta} = \frac{55.56 - 10 \times 6.545 \times 0.808}{437.611 - 10 \times (6.545)^2} = 0.291$$

$$\hat{\alpha} = 0.808 - 0.291 \times 6.545 = -1.097$$

$$s_f^2 = \frac{1}{9} [7.471 - 10(0.808)^2] = 0.1047$$

Probability Methods in Civil Engineering
IIT Kharagpur
Dr. Rajib Halty, IIT Kharagpur
15

What is their mean of this x' which is 6.545? \bar{y} is 0.808. So, this estimate of this beta hat as if you refer to this expression of this least square estimates then we can get that estimate of this beta hat will be 0.291 and the alpha hat will be minus point oh sorry minus 1.097 and the variance of this y which is unconditional which is your 0.1047 now using this alpha hat and beta hat. What we can calculate? That these modeled values, we can get for this y_i . So, this alpha and beta whatever we have estimated now and now we will use this inputs this x_i' as this input as this x and we will get what is the model estimate of this y_i .

So, that we will get and after we get this one then we can calculate what is their square errors. So, this 0.51 minus 0.563 square will give you this value similarly for all 10 observations which you have been calculated and sum them up to get that sum of square error which is 0.164.


(Refer Slide Time: 39:07)

Example...Contd.

$$s_{Y|x}^2 = \frac{0.164}{10-2} = 0.0205$$
$$s_{Y|x} = \sqrt{0.0205} = 0.143$$
$$r^2 = 1 - \frac{0.0205}{0.1047} = 0.804$$

□ The mean value function and standard deviation is:

$$E(Y|x) = -1.097 + 0.291 \ln x$$
$$s_{Y|x} = 0.143$$

 Probability Methods in
Civil Engineering
IIT KGP

Dr. Rajib Halty, IIT Kharagpur

16

Now, using that 0.164 we know that this is conditional variants of that target variable Y is that sum of square error divided by n minus 2 and here you know that two means that in number of input variable is one. So, 1 plus 1 it is 2 that is for just now we have seen for this multiple linear regression that it is that sum of square error divided by n minus m minus 1. So, m plus 1 degree of freedom is lost; here this 10 minus 2. So, if we calculate this one. So, these conditional variants will be 0.0205. So, this conditional standard deviation positive square root of that which is 0.143 and that percentage of that reduction of this variants. So, this 1 minus what you got for this conditional 0.0205 and what was the unconditional which is 0.1047 which is equals to 0.804.


Finally, the equation - that final equation that we get the mean value mean value function and the standard deviation is that expectation of Y given x equals to minus 1.097 which is the estimate for this alpha plus 0.291 estimate for the beta into log natural of that x. And, conditional variants given the x is constant, sorry, conditional standard deviation given that x is constant which is equals to 0.143. As we have seen and this square that you got is that 0.804. So, you can say that there that 80.4 percent of this variability has been explained through this model.

(Refer Slide Time: 40:54)

Correlation

□ Definition

- Correlation is a statistical technique that can show whether and how strongly pairs of variables are related.
- The study of the degree of linear association between two random variables is called correlation analysis.
- The accuracy of a linear prediction will depend on the correlation between the variables.

 Probability Methods in
Civil Engineering: H34.4
H315

Dr. Rajib Maiti, IIT Kharagpur

17

Well, now, we will just take that correlation part that we have discussed earlier and you know that the correlation is a statistical technique that can show whether and how strongly the pairs of the variables are related. The study of the degree of linear association between two random variables is called this correlation analysis and the accuracy of a linear prediction will depend on the correlation between the variables.

Now, we can in this regression context, what we can say is that, if we say that what we have modeled and what we have observed is this two are linearly associated and that linear association is stronger enough then, we can say that yes that just now we have seen that in terms of this percentage reduction of these variants in through this r square which is we can say that. So, that much variability is can be explained through that developed model. So, like that this accuracy of the linear prediction will depend on this correlation between those variables.


(Refer Slide Time: 42:03)

Correlation

□ In a two-dimensional plot, the degree of correlation between the values on the two axes is quantified by the so called correlation coefficient, which is given by :

$$\rho = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y} = \frac{E[(X - \mu_x)(Y - \mu_y)]}{\sigma_x \sigma_y}$$

where, E is the expected value operator,
and Cov means covariance

 Probability Methods in
Civil Engineering: H34.4
H34.5

Dr. Rajib Maiti, IIT Kharagpur

18

In a 2 dimensional plot now, if we just take that only one input here in case of the simple linear regression the degree of correlation between the values on two axis is quantified by so called correlation coefficient which is given by this equation. We have discussed earlier that is correlation coefficient is the covariance between x and y divided by standard deviation of x multiplied by standard deviation of y and the co-variants of x. You know that is x minus the is a expectation of x minus μ_x into y minus μ_y ; where this e is the expected value of this operator and C o v is that operator, means the covariance between these two; basically the sum. In that, after we develop this regression model we have to see that what is the correlation coefficient between that what we have observed Y and what we have modeled through.

So, basically here even though we are expressing this one, just to relate our earlier discussion x and y basically, we have to see it for this y i and that y estimates that is the y hat from that regression expression.


(Refer Slide Time: 43:11)

Correlation

□ The correlation coefficient may also be estimated by :

$$\hat{\rho} = \frac{1}{n-1} \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_x s_y} = \frac{1}{n-1} \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{s_x s_y}$$

where \bar{x}, \bar{y}, s_x , and s_y are respectively the sample means and standard deviations of X and Y.

 Probability Methods in Civil Engineering: H3L4
Dr. Rajib Maiti, IIT Kharagpur 19

Now, this correlation coefficient may also be estimated by this rho hat; now this estimation that we are getting 1 by n minus 1 i equals to 1 to n x i minus x bar y i minus y bar divided by S x S y. So, this one you can see that this expression can be written as this i equals to 1 to n x i into x i y i minus n into x bar into y bar. These are the means of this x and y whatever we have observed in this data divided by that S x into S y where these things x bar y bar Sx and Sy are the sample means and standard deviation of x and y respectively.

(Refer Slide Time: 43:57)


Correlation

□ We have:

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

□ Using the above two equations, we can rewrite the equation of correlation coefficient as:

$$\hat{\rho}_{\hat{\beta}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \frac{s_x}{s_y} = \hat{\beta} \frac{s_x}{s_y}$$

 Probability Methods in Civil Engineering: H3L5
Dr. Rajib Maiti, IIT Kharagpur 20

So, we have seen also that this in the simple linear regression that the estimate of this beta is equal; is having this form of this. If you just put this one in whatever in the expression of this rho then, we can get this expression; that this rho hat is equals to that beta hat multiplied by this ratio of this Sx and S y. So, that ratio between the standard deviation of x and standard deviation of x and standard deviation of y multiplied by this beta hat will give you that estimate of this correlation coefficient.

(Refer Slide Time: 44:39)

Correlation

□ Also we have

$$S_{y|x}^2 = \frac{1}{n-2} \left[\sum_{i=1}^n (y_i - \bar{y})^2 - \hat{\beta}^2 \sum_{i=1}^n (x_i - \bar{x})^2 \right]$$

□ Further, by substituting value of β in the above relation of co-variance, we can write :

$$\begin{aligned} \hat{\text{Var}}(Y|x) &= \frac{1}{n-2} \left[\sum_{i=1}^n (y_i - \bar{y})^2 - \hat{\rho}^2 \frac{s_y^2}{s_x^2} \sum_{i=1}^n (x_i - \bar{x})^2 \right] \\ &= \frac{n-1}{n-2} s_y^2 (1 - \hat{\rho}^2) \end{aligned}$$

Probability Methods in Civil Engineering
Dr. Rajib Halty, IIT Kharagpur

Now, also what we have seen that the conditional variants of y given that x is equals to 1 by n minus 2 that multiplied by i equals to 1 to n y i minus y barsquare minus beta hat square i equals to 1 to n into x i minus x y bar square.

This expression that we have seen in that simple linear regression to express what is the conditional variant. So, if we just put that estimate of this beta hat in terms of their correlation coefficient and their variance of those y and thus and x then this conditional variants that we can write that 1 by n minus 2 this expression and in place of this beta hat square, we can write that this rho hat, sorry, this will be rho hat the estimate of this correlation coefficient that rho square multiplied by this S y square and S x square.

These are the variants of y and this is the variants of x. So, if we just express this one then we can write that this is n minus 1 divided by n minus 2 s y square into 1 minus rho hatsquare. So, this one this part that we can take and we can express that what is coming

is that S_y^2 will come. So, this S_x and this one we can relate to this S_x^2 square variants that you know that this is this divided by $n - 1$ will give you the estimate of this S_x^2 square. So, that can be absorbed here and final expression that we are getting is that conditional variants of y given this x will be equals to $n - 1$ by $n - 2$ into S_y^2 into $1 - \rho^2$ square.

Now, if we can say that this estimate if we say that this equals to 1, sorry, if we say that this n is very large. So, that means, that when we say that there are large numbers of observation is available. So, for n when n is very large then we can say that this quantity can be equated to unity and that we can say that this will be equals to your S_y^2 into $1 - \rho^2$ square.

(Refer Slide Time: 47:39)

Correlation

from which we can write :

$$r^2 = 1 - \frac{n-2}{n-1} \frac{s_y^2}{s_x^2}$$

which can be approximated to r^2 for large n .

NPTEL Probability Methods in Civil Engineering IIT Kharagpur Dr. Rajib Halty, IIT Kharagpur 22

This is the final expression that we get which can be approximated. So, this one expression for this n equals to large or for any n if we just consider this factor to be multiplied that which can be approximated to this r^2 square for this large n .

What we have seen earlier that this r^2 square that is we have explained in terms of this percentage reduction of these variants which is equals to $1 - \rho^2$ the conditional variants divided by conditional divided by unconditional variants that was the square percentage of the reduction or that there is a percentage. How much is explained through that regression model. So, here we can see that for if the n is large than this quantity can be approximated to this r^2 square.

(Refer Slide Time: 48:29)


Example

Q. From the following results, obtain the two regression equations and estimate

- The yield of crop, when the rainfall is 22 cm
- The rainfall, when the yield is 600 kg.

	Mean	Standard Deviation
Yield (kg)	408.4	31.8
Rainfall (cm)	24.7	4.5

Co-efficient of correlation between yield and rainfall = 0.54.

 Probability Methods in Civil Engineering: H3R4
Dr. Rajib Maiti, IIT Kharagpur 23

So, what we have seen so far is that, if we know that the correlation coefficient between the variables that we are modeling that is what is your target variable and what is your input variable if you know, that what is their correlation coefficient between then and if we know that, what is their respective mean then, basically what we can do? We can develop, we can get the estimate of those parameters of the regression.

So, we can basically develop their regression equation. One such example we will just see now. This example is on that; from the following results, obtain the two regression equations and estimate the yield of crop when the rainfall is 22 centimeter and estimate the rainfall when the yield is 600 kg.

Basically, this is the relation between this yield of the irrigation and what is the rainfall in terms of this depth of the rainfall in centimeter and the mean of this yield? The mean yield is 408.4 kg and this mean rainfall is 24.7 centimeter.

So, the standard deviation of this yield is 31.8 kg and standard deviation of this rainfall is 4.5 centimeter and the correlation coefficient between the yield and rainfall is point 54. So, this information is available to us. So, if you know this one then we have to develop the two regression equations, that is the what is the expression for this yield given the rainfall and what is the expression for the rainfall given by yield. So, this both this equation we can use with the help of this information of this correlation coefficient.

(Refer Slide Time: 50:23)

Example...Contd.

Sol.:

Let Y be yield and X be rainfall.


So, for estimating the yield, we have to run the regression of Y on X and for the purpose of estimating the rainfall, we have to use the regression of X on Y.

Given

$$\bar{X} = 24.7, \bar{Y} = 408.4, \sigma_x = 4.5, \sigma_y = 31.8, \rho = 0.54$$

Therefore, regression coefficients are :

$$\beta_{YX} = 0.54 \cdot \frac{4.5}{31.8} = 0.0764 \quad \beta_{XY} = 0.54 \cdot \frac{31.8}{4.5} = 3.816$$

 NPTEL Probability Methods in Civil Engineering IITLS

Dr. Rajib Halty, IIT Kharagpur

24


Let that Y be yield and X be the rainfall. So, for estimating the yield we have to run the regression of Y on X and for the purpose of the estimating the rainfall we have to use the regression of X on Y and the information that you know is that mean of this X mean of this Y standard deviation of X standard deviation of Y and their correlation coefficient between them.

So, this beta X given Y that is when we are when we regress X on Y. So, this is that correlation coefficient divided by their ratio of their standard deviation. So, what we get that 0.0764 and on the other hand when we regress that Y on X then that beta coefficient will be 3.816.

(Refer Slide Time: 51:17)

Example...Contd.

- Hence, the regression equation of Y on X is
$$y - \bar{y} = \beta_{YX}(x - \bar{x})$$
$$Y - 408.4 = 3.816(X - 24.7)$$
$$Y = 3.816X + 314.145$$
- Similarly, the regression equation of X on Y is
$$x - \bar{x} = \beta_{XY}(y - \bar{y})$$
$$X - 24.7 = 0.0764(Y - 408.4)$$
$$X = 0.0764Y - 6.502$$

 NPTEL Probability Methods in Civil Engineering IIT Kharagpur


Dr. Rajib Halty, IIT Kharagpur 25

So, the regression equation for this Y on X will be $y - \bar{y} = \beta_{YX}(x - \bar{x})$. So, the coefficient when we are regressing Y on X that is denoted like this that β_{YX} multiplied by $x - \bar{x}$. So, that we have estimated 3.816. So, after rearranging this we can get the equation like $y = 3.816x + 314.145$. So, this is the regression equation for the Y on X similarly we can regress the equation of this X on Y which finally, we get that $x = 0.0764y - 6.502$.

(Refer Slide Time: 52:26)

Example...Contd.

- So, when $X = 22$,
$$Y = 3.816 \times 22 + 314.145 = 398.1$$
- And when $Y = 600$,
$$X = 0.0764 \times 600 - 6.502 = 39.34$$
- Hence, the estimated yield of crop is 398.1 kg and the estimated rainfall is 39.34 cm.

 NPTEL Probability Methods in Civil Engineering IIT Kharagpur

Dr. Rajib Halty, IIT Kharagpur 26

With this one if we get, if we use this expression and then the question was given that when x is equals to 22, what is the yield? That is when the rainfall equals to 22 centimeter what is the yield? So, putting that in this expression we get that Y is equals to your 398.1 and when that yield is equals to 600 then the estimate of this rainfall is your 39.34. So, the estimated yield of this crop is 398.1 k g and the estimated rainfall is 39.34 centimeter. So, this is when this yield is when rainfall is 22 centimeter and this rainfall is estimated like this when the yield is 600 kg.

So, in this lecture or including this last lecture we have discussed the regression different regression technique including their simple linear regression, multiple linear regressions then non-linear regression and that in terms of this correlation. How we can estimate that? We have discussed. So, in this entire module what we whatever we have seen in this probability and style statistics. We have started with the sample statistics then we have covered this hypothesis testing. Then, how to test what the data follows? What distribution through this probability paper and different test different statistical test? To test that, what is the distribution of this of the parameter and finally, we have covered the regression analysis and correlation, thank you.