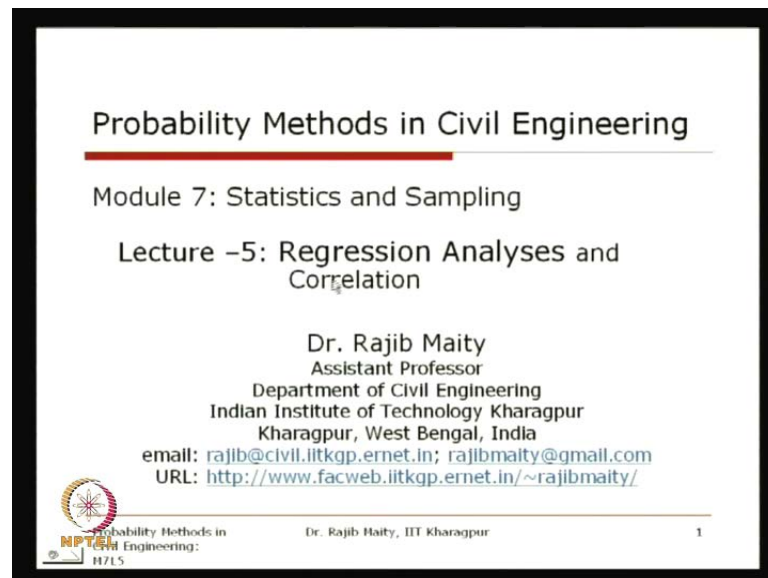


Probability Methods in Civil Engineering
Prof. Dr. Rajib Maity
Department of Civil Engineering
Indian Institute of Technology, Kharagpur

Lecture No. # 39
Regression Analyses and Correlation

Hello and welcome to this lecture. In this lecture, may be in this or the next lecture, we will cover the topic on this regression analyses.

(Refer Slide Time: 01:55)




Probability Methods in Civil Engineering

Module 7: Statistics and Sampling

Lecture -5: Regression Analyses and Correlation

Dr. Rajib Maity
Assistant Professor
Department of Civil Engineering
Indian Institute of Technology Kharagpur
Kharagpur, West Bengal, India
email: rajib@civil.iitkgp.ernet.in; rajibmaity@gmail.com
URL: <http://www.facweb.iitkgp.ernet.in/~rajibmaity/>

 Probability Methods in Civil Engineering: H7L5

Dr. Rajib Maity, IIT Kharagpur 1

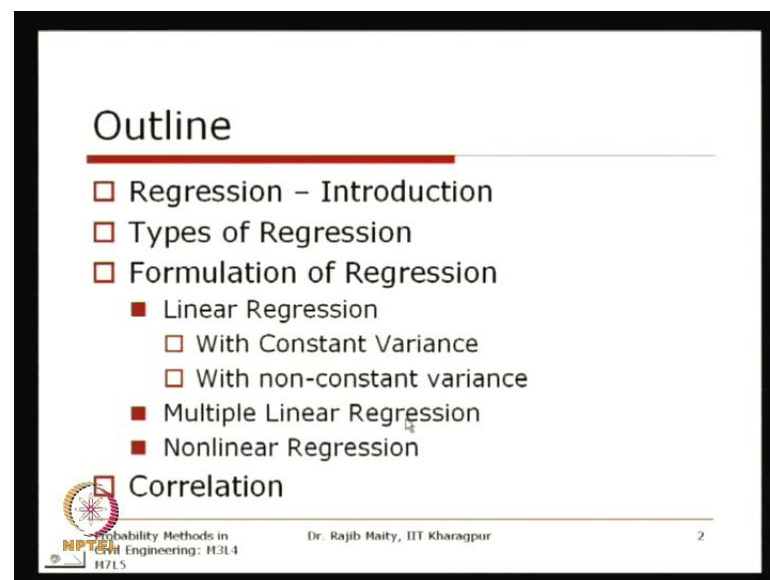
There are different types of regression analyses say, first we will start with that simple linear regression and after that there are different types. Multiple regressions are there and then non-linear non-linear regressions are there. So, we will see and basic fundamental things; basic concept we will understand. Basically, when in a, in many application fields, obviously, including civil engineering there are many random variables are there which are supposed to have some relationship to their and in this, through this analyses we tried to capture that we try to model that relationship.

Now, if the relationship is linear then we generally go for this linear regression and sometimes we have seen that maybe the linear relationship is not sufficient. So, there we

have to go to the non-linear regression analyses. Sometimes, the target variable is dependent only on one variable or sometimes that response variable or the target variable can depend on more than one independent variable. So, in that case we generally go for these multiple regressions. So, all these things we will learn in this lecture or this may continue to the next lecture also. So, this is our today's lecture title is regression analyses and correlation and this correlation means here that we have already discussed earlier that this correlation when we discuss this random variable and all.

So, here also we will see that how this regression analyses. In this regression analyses correlation is an important part. So, we will just see in the light of this regression analyses also towards the end of this lecture.

(Refer Slide Time: 02:24)



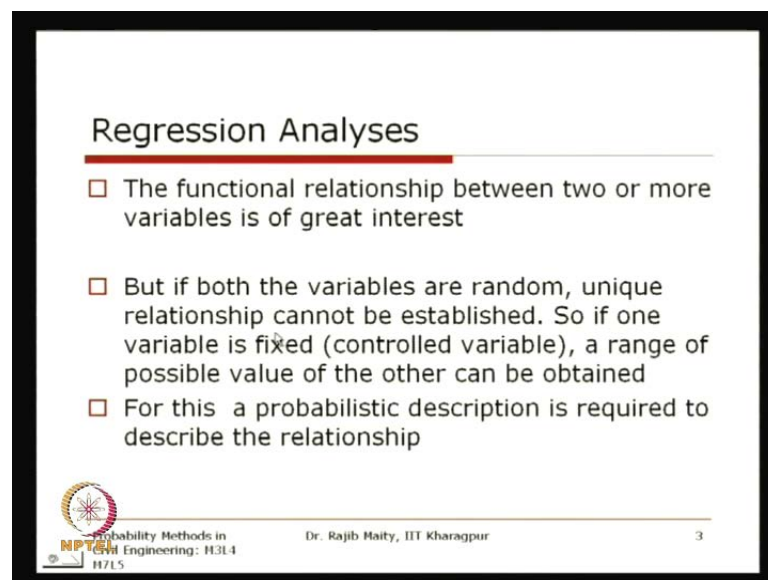
So, our outline of this today's lecture is that first we will go through some introduction and then we will discuss about some that different types of regression then you formulation of this regression in this there are linear regression as I mentioned and in this linear regression also it may have the constant variance or it may have the non constant variance. So, this non constant may be the variable variance may be the other word, but, just to avoid to similar word. So, it has used as this non constant variance.

So, this non constant variance and this constant variance means in general for the linear regression when we refer to we refer to this constant variance we means over the entire range of the dependent variable the variance of the response variable remains same that is

what is the I can say that by default case, but, sometimes or it can be observed that these variance may also vary over the different range of that dependent variable. So, in that case, we have to go for this non constant variance also.


Then, if there are more than one dependent variable then we have to go for this multiple linear regression and if the relationship we see that may not be linear sometimes some other non-linear relationship may have better, can better extend the target variable then we can go for this non-linear regression and then as I told that there is... So, we will see that correlation basically this will be a major that how strong the relationship has been captured through that model that we have developed through this regression; so, that we will see.

(Refer Slide Time: 04:09)



Regression Analyses

- ☐ The functional relationship between two or more variables is of great interest
- ☐ But if both the variables are random, unique relationship cannot be established. So if one variable is fixed (controlled variable), a range of possible value of the other can be obtained
- ☐ For this a probabilistic description is required to describe the relationship

 Probability Methods in
Civil Engineering: H3L4
H7L5

Dr. Rajib Maity, IIT Kharagpur

3

Well, in this regression analyses the fundamental that, sorry, the functional relationship between two or more variables is of great interest as I mentioned that there may be many variables which are which we can see that there could be a relationship by the linear or non-linear and this kind of relationship basically if we just take the observed data and plotted through some scatter diagram and then itself by visual inspection itself we can see that there are whether there are some types of relationship is there or not. So, if we can see then we can think of this type of regression analyses to capture that particular relationship.

So, here if both the variables are random, a unique relationship cannot be established. So, you know that a unique relationship here what is meant is that it may not be that one-to-one relationship; there could be some event where there could be some randomness in both the variables. So, if one variable is fixed and that is known, that is termed as a control variable or that what I mentioned is that is the dependent variable; the range of possible values of the other can be obtained through this analysis. For this, a probabilistic description is required to describe this relationship and basically this is what you will get through this regression analysis.

(Refer Slide Time: 05:45)

Regression Analyses

□ Relationships between variables :

- How does the strength of a material depend on temperature?
- How does the compressive strength of concrete depend on water cement ratio?
- How strong is the link between rainfall and runoff?

NPTEL Probability Methods in Engineering: H7L5 Dr. Rajib Haity, IIT Kharagpur 4

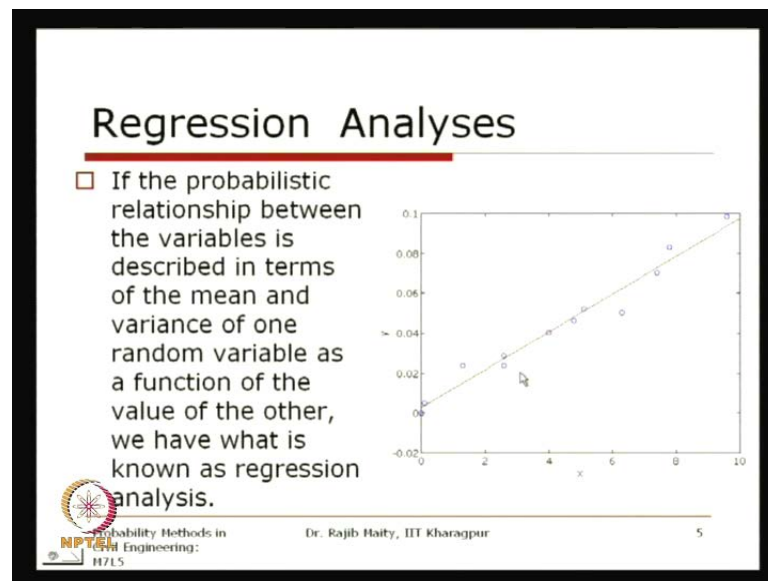
So, in this, the type of question particularly if I concentrate to this different field of application in civil engineering then this type of analysis will give me the answers to a kind of this type of question say that how does the strength of material depend on the temperature. So, if the temperature I vary. So, how the strength of material whether it will increase or decrease or how the relationship is.

Second is say that how does the compressive strength of the concrete depend on the water cement ratio. So, if I increase the water cement ratio then what will happen to the compressive strength or if I decrease it what will happen. So, these are some two variables are considered. Similarly, what we can say that whether that target variable here is the compressive strength may have instead of this only that water cement ratio. There could have been other factors as well that can be influencing to this. So, then

what will happen? That one target variable and more than one dependent variable - so that, multiple regression can come into the picture.

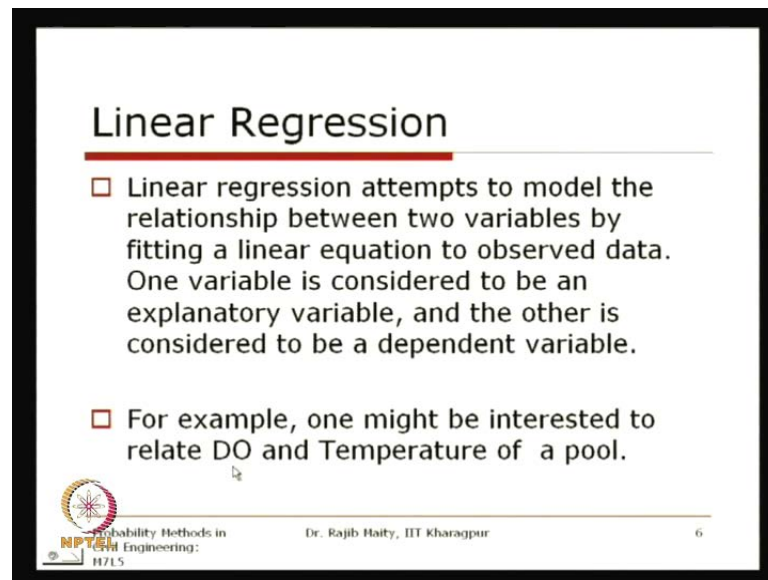
And, third: so that, how strong the link between the rainfall runoff for a given catchment or for a given area. So, how. So, rainfall and runoff if the rainfall is more runoff can be more. So, how strong is that relationship? So, this type of answer we can get through this regression analysis.

(Refer Slide Time: 07:16)




The probabilistic relationship between the variable is described in terms of the mean and variance of one random variable as a function of the value of the other we have what is known as the regression analysis. So, say for example, as I was telling just by if I just plot that through a scatter plot the what is the observed data that we are having the paired observed data paired in the sense here that we are talking about the two variables first. So, this is one variable is x and other one is the y . Now, if I just plot it, these blue circles, you can see that this is the paired data and. So, we can see that if x increases y can also increase and vice versa if x is decreasing y is decreasing. So, whether now can we just estimate one relationship between this x and y . So, that estimate that estimate of this functional relationship is that regression analysis that we will get through this analysis.

(Refer Slide Time: 08:24)



Linear Regression

- Linear regression attempts to model the relationship between two variables by fitting a linear equation to observed data. One variable is considered to be an explanatory variable, and the other is considered to be a dependent variable.
- For example, one might be interested to relate DO and Temperature of a pool.

 Probability Methods in Civil Engineering: H7/L5

Dr. Rajib Haity, IIT Kharagpur

6

So, first we will take that linear regression for examples are the diagram that is shown here we can see and we can expect that there could be a linear relationship can have it here. So, but, in many other cases where if just is looking this scatterplot we can see that initially it may be increasing and later on it may not increase in that rate. So, there could be we can expect that there might be a non-linear relationship can happen. So, the first what we are taking up is that linear regression; where the expectation that the relationship is linear between the dependent variables and the target variable.

So, the linear regression attempts to model the relationship between two variables by fitting a linear equation to the observed data one variable is considered to be an explanatory variable and other is considered to be a dependent variable. So, that is what our target. So, in this example that we have seen what we can use is that this variable x we can use as to be that your dependent variable and this is a y is my target variable. So, I can use the information of x and I can model this y it can be it could not be opposite also if we can we if we estimate a x with respect to the variable y . So, then we generally say that that x is regressed on y and in other way the y is regressed on x .

For example, that one might be interested to relate the dissolved oxygen and the temperature of a pool. So, whether the dissolved oxygen and temperature these two data is generally first collected and then we can see that whether their relationship how the relations how they vary with respect to each other whether they must whether in the sense

that I can see it in both sides whether that DO given the temperature or the temperature given what is the DO, but, sometimes in case of this the practical consideration may be we are interested to know that our what is our target what is the what should be the dependent variable and what should that target variable for example, the example that is given here the dissolve oxygen and the temperature generally what we see is that temperature we use as a dependent variable and this dissolve oxygen is the target variable. So, this depends on the in what area in what practical field that we are that we are applying this analyses.

(Refer Slide Time: 11:10)

Basic Formulation of Linear Regression

- Regression with constant variance
 - Let us consider a pair of data X, Y plotted on a scatter diagram.
 - From the figure it can be noted that the possible values of variable Y depends on the other variable X
 - So to analyze the data for Y (calculating variance and mean), we take into consideration the change in X .
 - Also we can see that there is a general tendency for values of Y to increase with X

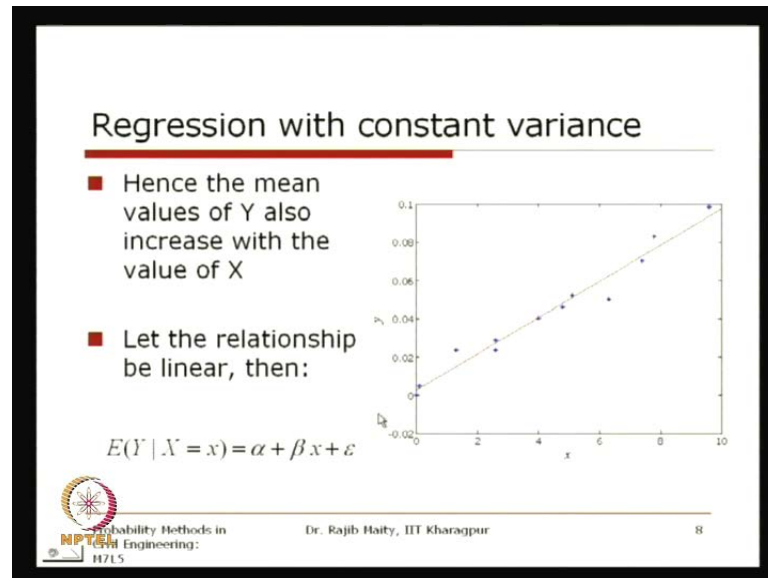
NPTEL Probability Methods in Engineering: H7L5 Dr. Rajib Hait, IIT Kharagpur 7

The basic formulation of linear regression with the constant variance is first. So, here as I was stating as a starting that when we are taking that the constant the variance of the dependent variable over the entire range of the dependent variable it remains constant. So, in that case we generally say that this regression with a constant variance and by default when you say the regression analyses we generally mention that it is with the constant variance. So, the non constant case is a special case that we will take that we will see after some time.

So, in case of this regression with constant variance let us consider a pair of data XY plotted on the scatter diagram as just few slides before you have seen from the figure it can be noted that the possible values of the variable Y depends on the other variable X . So, to analyze the data for Y calculating variance and mean we take into consideration the

change in X and also we can see that there is a general tendency of the values of Y to increase with the X. So, these are some of this example is given with respect to that plot the scatter plot that we have seen few slides before.

(Refer Slide Time: 12:35)



So, here again the similar plot has been shown here. So, here that one variable is X and other one is the Y. So, here we are taking the case that we will regress Y on X. So, X is our dependent variable and Y is our target variable. So, here you can see that when X increases Y also increases and vice versa. So, we have to fit a linear relationship between this X and Y. So, hence the mean value of mean values of Y also increases with the value of the X. So, as X increases that mean value or the in the statistical sense the expected value of the Y also increases. So, the relationship let the relationship be linear because we are discussing this linear regression now. So, the expected value of the Y given a particular value X. So, you know. So, this is the conditional expectation. So, if I just take what is the expectation of the Y you know the expectation of the y means without any other information. So, whatever the Y we see that it can from this diagram we can see that it varies from 0 to 0.1 say. So, whatever the values the range that we see we will just take its mean and that is the expected value of the Y.

Now, this when you are fitting this relationship; that means, ; that means, it is a condition on the given value of this X now if I give some value of this X at 6. So, a in this part what is the expected value of this Y. So, this is now becomes the condition and this


conditional mean is expressed through this linear relationship which is $\alpha + \beta x + \epsilon$. So, this is what you know this is the equation of that of the straight line plus some error term should be there to express that what is that value of that the mean value of this y .

(Refer Slide Time: 14:42)

Regression with constant variance

- α and β are constants
- ϵ is the possible error of measurement
- Variance of Y may be independent or a function of x
- This is known as linear regression of Y on X
- Now we have to estimate the parameters α and β of the regression line, such that it provides the best fit of the data

i.e if we have n paired of observation (x_i, y_i) , then determine \hat{Y} such that difference between y_i and \hat{y}_i is minimum



Probability Methods in
Engineering:
M7L5

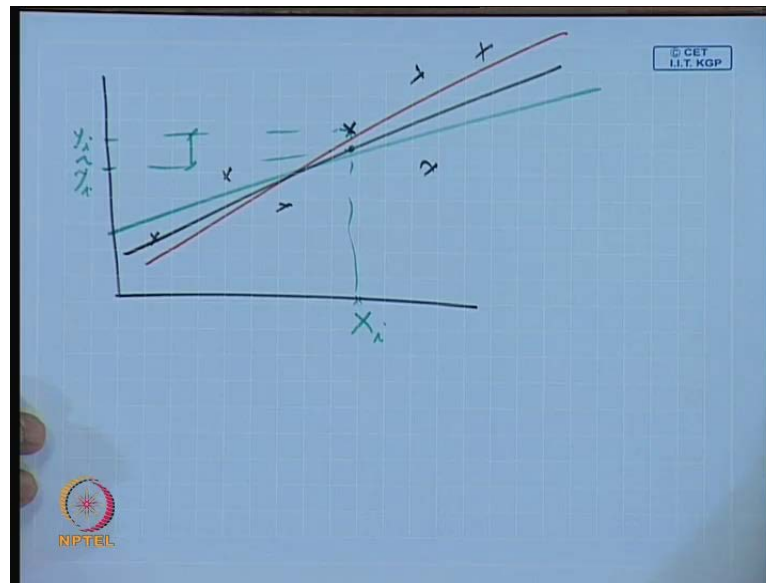
Dr. Rajib Haity, IIT Kharagpur

9

Now this α and β are the constant and ϵ is the positive error of this measurement of the sorry possible error of measurement. So, if when we take that data observe that data. So, there could be some in that measurement. So, there could be some errors. So, that error is expressed through this ϵ variance of Y may be independent or a function of X this is known as the linear regression of Y on X that is what I was telling. So, it is Y on X it can be expressed in other way also that is X on Y . So, the relationship will change that expectation of X given Y is equals to some constant plus the β multiplied by that your Y plus ϵ . So, that is the observational error.

So, now we have to estimate the parameters of this α and β of the regression line such that it provides the best fit of the data now this best fit of the data means if I just see this one this scattered diagram. So, there could be the various possible lines that I cannot think through these points now which line should be the best fit line. So, what is meant here is this.

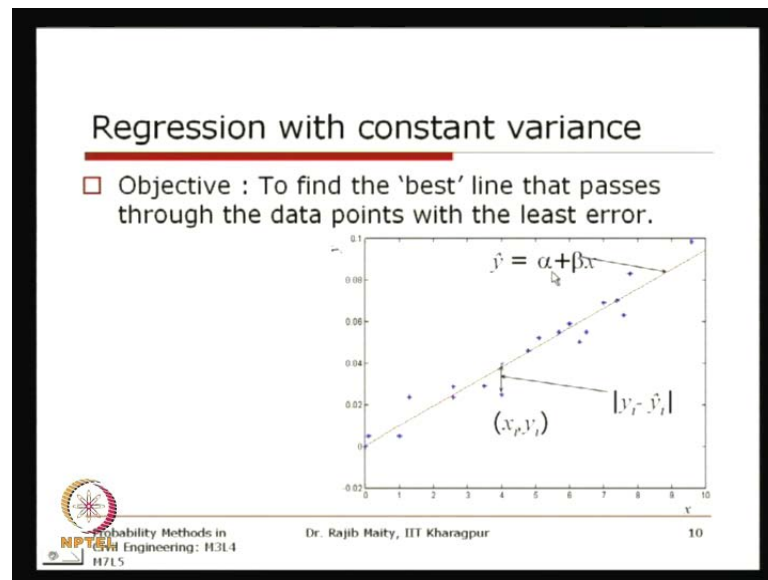
(Refer Slide Time: 16:09)



So,, if this is these are the data points then there could be the there could be some lines which canbedescribed throughdifferentstraight lines now out of these lines the possible lines which one should be the best fit now this to get that best fit. So, to get that best fit we have to follow somemethodology which is known as the method of least square to estimate thatto get that the line that is best fitting through this points and based on that we will get what is thethat estimate of those regression constants. So, this is what is mentioned here that issuch that its provide the best fit of the data that is if we have npaired of thisobservation $x_i y_i$. So, these are the paired observation and these are one pair is one point on that diagram thendetermine y_i cap in such that the difference between that y_i and y_i cap is minimum.

Now, what is this y_i and y_i cap is if I referred to this diagram is this. So, this is your that point where you can see that this is your somethis is that y_i and if the whatever suppose the this black line is your best fitline then this with respect to this x_i with respect to this x_i . So, the estimate is this one. So, this is your y_i cap. So, the difference between these two is the error which should be minimized. So, now, as close as thispoint to this observation and this is for all the points then that line should be the best fit line. So, the difference between that is why the difference between this y_i and y_i cap should be minimum to declare that the line is a best fitting through the data.

(Refer Slide Time: 18:23)



So, it is now explained here. So, our objective here to find the best line that passes through the data points with the least error. So, now, this blue stars are the observed data and this is the estimate of this regression line which is alpha plus beta x now this difference from what is the point that you can see and what is this corresponding point on this regression line there is a red line shown here is your error. So, this mode of this y_i minus \hat{y}_i is the absolute error for the point x_i, y_i . So, y_i is known is the observed one and that \hat{y}_i is what we will get that will be your \hat{y}_i .

(Refer Slide Time: 19:23)

Regression with constant variance

■ The constants α and β are found by minimizing the sum of the squared errors (Principle of Least Squares)

$$\Delta^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

$$\frac{\partial \Delta^2}{\partial \alpha} = \sum_{i=1}^n 2(y_i - \alpha - \beta x_i)(-1) = 0$$

$$\frac{\partial \Delta^2}{\partial \beta} = \sum_{i=1}^n 2(y_i - \alpha - \beta x_i)(-x_i) = 0$$

Dr. Rajib Haity, IIT Kharagpur

Probability Methods in Civil Engineering: H3L4 H7L5

So, the constants alpha and beta are found by the minimizing the sum of squared errors sum of squared errors and this is known as this principle of least square. So, what is done is that this is the error that is y_i minus \hat{y}_i this is the error and that error is squared and summed up for all the observations. So, in this diagram if we see, this is the error y_i minus \hat{y}_i and this is obtained for all these data points and this error for individual point is first square up and that. So, that square error is summed of for all the observations that is available. So, this is giving is the sum of square errors.

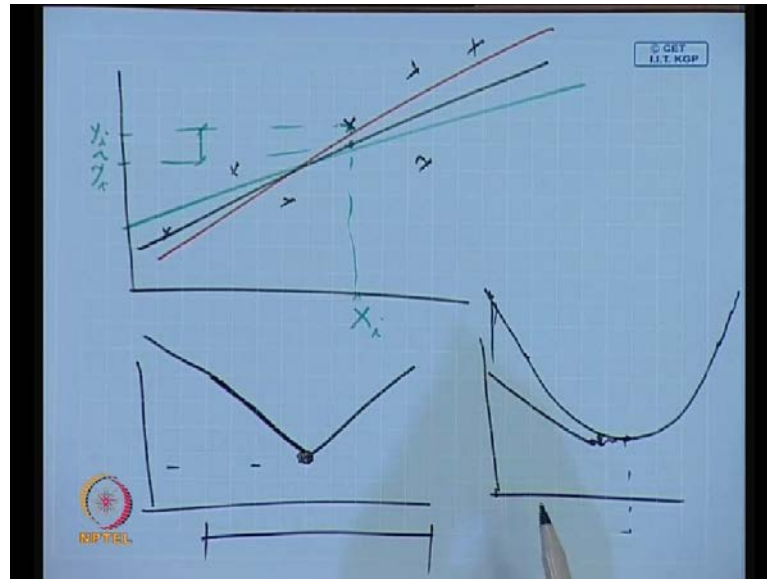
Ah now this sum of square error now if I just replaced this \hat{y} from that regression line which is your $\alpha + \beta x_i$ and we are taking this minus. So, y_i minus α minus βx_i whole square is give you that sum of square errors now to get that estimate of this alpha and beta. So, this error should be for this alpha and beta the value of alpha and beta should be such that because these are the two constants which is basically which is determined everything about that straight line. So, error - this quantity should be minimum.

Now, to ensure that we have to take this partial derivative of this sum of square error with respect to each this parameters alpha and beta and this has to be equated to 0. So, we are having two unknowns and we are having two simultaneous linear equation that we can solved to get what is the estimate of this beta before I proceed I need to take some time to explain this one why we have taken this square and this and we are using this as that total error.

So, because you can see the first thing the first direct thing that you can that you can have from this diagram is that. So, for some points the error will be negative and some point the error will be positive - the error will be positive depending on whether the point is below the regression line or above the regression line.

Now, when we are taking this square obviously, those sign is going because we are interested to this what is the deviation from this regression line whether it is on the positive side or on the negative side that we are not interested when we are looking for the best fit line. So, whatever the error that we get if we take the square; obviously, that sign will go, but, this can also be, think of that if we just take that absolute value of that error as it is shown it here then also that sign can go and if we just add them up then, what we will get is that also it will give an the absolute error summation of the absolute error.

(Refer Slide Time: 22:59)



But, generally when we go for this least square technique we take it to be the square and then we do this partial derivative this is because you know when we take the error now if we just see that error and that error if we take it as a linear function basically what we are if we are minimizing it then this one basically our point is suppose this is our target point now this is the over that the possible range of the parameters now when we take this absolute error then the change with respect to that parameter it will be the linear one and when it take it to this to the square or the. So, this will become basically a quadratic function.

Now, what will happen if if our estimates are far away from what is the optimum value. Then you know from the optimization technique. So, if it is far away then the next step basically it will go very close to that optimum value and once it comes to the optimum value then, the steps will be smaller steps. But, in this case generally that the step size are always same because this variation is linear here, but, means this is basically when you go for this optimization optimizing the parameters that time it has been seen that this taking square is better than this taking this linear function

So, that is why so far as that sum of it is we are all generally interested for this sum of square error. So, what is done in this principle of least square technique?

(Refer Slide Time: 24:50)

Regression with constant variance

So, we obtain the least-squares estimates of α and β as :

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

Dr. Rajib Maity, IIT Kharagpur

12

Well, we got this after, through this partial derivative we get these two simultaneous linear equations whereby solving them we can get the estimates like this that alpha cap. This is now, this cap symbol is given; when we are referring to that it is the estimate. So, this alpha cap if we can solve it and we can. So, that it will be the y bar minus beta cap x bar. This y bar and x bar are the mean of this observed data and this beta cap is the estimate of this beta. Can be shown that it is the summation of this xi minus x bar multiplied by yi minus y bar multiplication of them sum it over the all n observation divided by xi minus x bar square sum it over this all n observation.

(Refer Slide Time: 25:34)

Regression with constant variance

So, the least-squares regression line is :

$$E(Y | x) = \hat{\alpha} + \hat{\beta}x$$

Similarly, we may also obtain the least squares regression of X on Y, i.e. $E(X | y)$, using the same procedure

Dr. Rajib Maity, IIT Kharagpur

13

So, these two are the estimate of this alpha and beta. So, the least square regression line is the expected value of this y given that x is equals to alpha cap plus beta cap x. Similarly, we may also obtain the least square regression of this x on y as I was mentioning that is that expected value of this x given y using the same procedure. But, here it will come as their dependent variable will be y. Obviously, that alpha and this beta the estimate of these regression parameters; obviously, will change through that if we follow that procedure whatever we have done.

(Refer Slide Time: 26:12)

Regression with constant variance

- Now the conditional variance $\text{Var}(Y|x)$, (the variance about the regression line) can be calculated as follows:
- Here, the conditional variance is assumed to be constant within the range of 'x' of interest. So, we have

$$s^2_{Y|x} = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

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

So, now the conditional variance - that is now the conditional variance of this now the variance of y given x. So, whatever we have got if that just now is that expected value of y given x. So, now, we are interested to know, what is the conditional variance of y given that x? So, the variance about the regression line basically, what we meant here is that if I just referred to this diagram that this is the, if this is the y this is y. So, we can see that it is varying from this 0 to 1. So, whatever the y observed data that we have got that we got and which we know how to obtain that it is a sample estimate of this variance; if we do so that will give you the variance of the y.

Now, after we get this regression line now, what is the variability of the y with respect to the regression line? So, basically we are looking through this access and we see that how it is varying across this regression line. So, that is what is referred to as means pictorially as this variance of y given x now if we want to estimate that one if we want to calculate that

one this can be calculated as follows. Here, the conditional variance is assumed to be constant within the range of this x . So, we have this $s^2_{y|x}$ equals to $\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$.

So, this is basically the estimate from this regression and there are $n-2$ to make it that what it is called that unbiased and this you know that in the standard deviation we have seen that one degree of freedom is lost and that is why we make it that $n-1$ that we discussed in the earlier lectures and here one more degree of freedom is lost when we are estimating that regression line.

Basically, there are if we see that there are two parameters that both α and β has to be estimate through this regression line and that is why the two degrees of freedom is lost. So, make this estimate unbiased it is $\frac{1}{n-2}$ that we have to make. So, we can just do this we can sometimes for this we can make that $y_i - \bar{y}$ that is a mean of y whole square minus β cap square $\sum_{i=1}^n (x_i - \bar{x})^2$. So, this is just from this equation and you can see that this is basically that error and so, sum of square errors. So, which is that $s^2_{y|x}$ that is that conditional variance of y given x .

(Refer Slide Time: 29:12)

Regression with constant variance

- Taking into account the general trend with X , the physical effect of the linear regression of Y on X can be measured by the reduction of the original variance of Y , s_y^2 and it is represented as :

$$r^2 = 1 - \frac{s^2_{y|x}}{s_y^2}$$

where, $s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$
is the sample variance of Y .

NPTEL Probability Methods in Mechanical Engineering: ME314 Dr. Rajib Maity, IIT Kharagpur 15

Now, taking into account the general trend with x the physical effect of this linear regression, what actually is happening through this linear regression is that, y on x that is the regression y on x can be measured by a reduction of the original variance of this y . So, the original variance of this y you know that which is that a s_y^2 square which I can

get the data of this y and we can estimate what is that what is this variance from the sample and that. But, through this regression when we do it that there is a regression there is a reduction in that variance. So, which can be expressed through this one minus that variance of this y given x divided by variance of this variance of y , sorry, this square will that $S_{\text{power square}}$. So, what you can see is that this is the original variance that was there in the data y and this is the variance after the regression that we that we got. So, this is basically how much is the reduction then that reduction is that what is the total minus what we got after this regression divided by what was their the total.

So, this is that reduction in that variance and later on, we will just show you that this can be approximated to basically the correlation coefficient. Obviously, the square root of this one is, can be approximately equal to the correlation coefficient and that is basically the measure of how strong the relationship that we have measured for this S_y square; means this is the one by that you know it is the sample estimate of this variance of the data y . So, $1/y$ by n minus 1 summation of from 1 to n y_i minus \bar{y} , sorry, \bar{y} square.

(Refer Slide Time: 31:07)

Example

Q. Given set of data where the shear strength in kPa obtained from sample taken from 10 different depths of clay stratum. Assume the variance is constant with depth and determine the mean and variance of the shear strength as a linear function of depth.

Depth x_i	Strength(kPa) y_i	Depth	Strength(kPa)
2	12	6	38
3	25	7	45
3	24	7	65
4	26	8	70
5	40	9	75

Dr. Rajib Maiti, IIT Kharagpur

We will take one problem. Whatever we have discussed through the, for this linear regression, given the data where the shear strength in kilo pascal obtained from the sample taken from 10 different depths of this clay stratum, assume that the variance is constant with the depth and determine the mean and variance of the shear strength as a linear function of the depth.

So, here you can see that there are depths are given the depth at 2,3 again this 3. So, there are 10 such depths are taken. There are some depths are same; you can see here and we are getting this data. So, for the strength that in kilo pascal, we are having these 10 different data set; this is the depth and this is basically this depth is going to 3,3,4,5,6, 7,7,8,9 and these are the corresponding strengths. So, this 10 data set that we are having and we will follow whatever we have discussed just to find out the relationship between strength and the depth.

So, we will regress the strength on depth. So, our we can say that our variable y is here the strength and the x is here depth.

(Refer Slide Time: 32:25)

Example...Contd.

Sol.:

- **Step 1:**
 - XY, X^2, Y^2 are determined as shown in table in next slide
- **Step 2:**
 - Then $\Sigma X, \Sigma Y, \Sigma XY, \Sigma X^2, \Sigma Y^2$ are also determined
- **Step 3:**
 - α and β are calculated using the formula discussed before

NPTEL Probability Methods in Civil Engineering: Dr. Rajib Halty, IIT Kharagpur 17

So, to get this estimate through this least square technique, we will first get the estimate of this parameter; that is XY, X^2, Y^2 are determined and this will be, we will show in this table and then the summation of this X summation of Y summation, XY summation of X^2 and Y^2 and then, the alpha beta are obtained for using the formula that we obtained through that from that least square estimate.

(Refer Slide Time: 32:52)

Example...Contd.

No:	x_i	y_i	$x_i y_i$	x_i^2	y_i^2	$\hat{y}_i = a + \hat{\beta} x_i$	$(y_i - \hat{y}_i)^2$
1	2	12	24	4	144	-0.515	156.624
2	3	25	75	9	625	-0.224	636.245
3	3	24	72	9	576	-0.224	586.797
4	4	26	104	16	676	0.067	672.513
5	5	40	200	25	1600	0.358	1571.472
6	6	38	228	36	1444	0.649	1395.078
7	7	45	315	49	2025	0.940	1941.257
8	7	65	455	49	4225	0.940	4103.645
9	8	70	560	64	4900	1.231	4729.127
10	9	75	675	81	5625	1.522	5398.957
Σ	54	420	2708	342	21840		21191.716

So, this is the data for this different depth and for this 10 data sets there are x_i, y_i square y_i square and then, these things we will just see. So, first we are having up to this and we are having their summation also.

(Refer Slide Time: 33:14)

Example...Contd.

$$\bar{x} = \frac{54}{10} = 5.4 \quad \bar{y} = \frac{420}{10} = 42$$
$$\hat{\beta} = \frac{2708 - 10 \times 5.4 \times 42}{342 - 10 \times (5.4)^2} = 8.73$$
$$\hat{\alpha} = 42 - 8.73 \times 5.4 = -5.143$$
$$s_f^2 = \frac{1}{9} [21840 - 10(42)^2] = 466.67$$

So, up to this of this table we know and using this information that is what is the power x bar is 5.4 y bar is 42, sorry, it is 42. This will be 42×42 by 10; so, it is 42.

So, this beta 1 you know that this expression we will use. So, and we will get that estimate of this 8.73 and alpha cap is the estimate of the minus 5.143 and this s_y^2 square there is a variance the or the total variance I can say now the total variance of this y is 466.67 now.

(Refer Slide Time: 33:49)

Example...Contd.

$$s_{Y|X}^2 = \frac{358.73}{10-2} = 44.8841$$

$$s_{Y|X} = \sqrt{44.8841} = 6.696$$

$$r^2 = 1 - \frac{44.841}{466.667} = 0.9039$$

□ Thus the mean value function and standard deviation is:

$$E(Y|X) = -5.143 + 8.730 \cdot x$$

$$s_{Y|X} = 6.696$$

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


So, this now S_y given that x is your this 44.88 and the standard deviation is 6.696 and this square is equals to 1 minus these how much is the reduction is that 0.9039; and this alpha and beta that you can see it here. So, alpha is minus 0.50.143 and beta is this 1. So, this regression equation comes like this; that expected value of this y even x is that minus 0.50.143 plus 8.730 x. Now, using these things basically this relationship in the table we got this expression first. So, we are putting this x input and that alpha estimate and beta estimate and we get this one.

From here we are getting what is their error of this that is y minus y caps. So, this minus this and that square will give you that basically what is this one that we get that error square now if we sum it up this is basically the sum of square error and we are using that information to estimate this what is that reduction in this that variability variance in y. So, this one we have seen. So, this is finally, that expected mean is expressed through this expression and expected variance is 6.697 and obviously, this is constant over the entire range of x.

(Refer Slide Time: 35:22)

Regression with non-constant variance

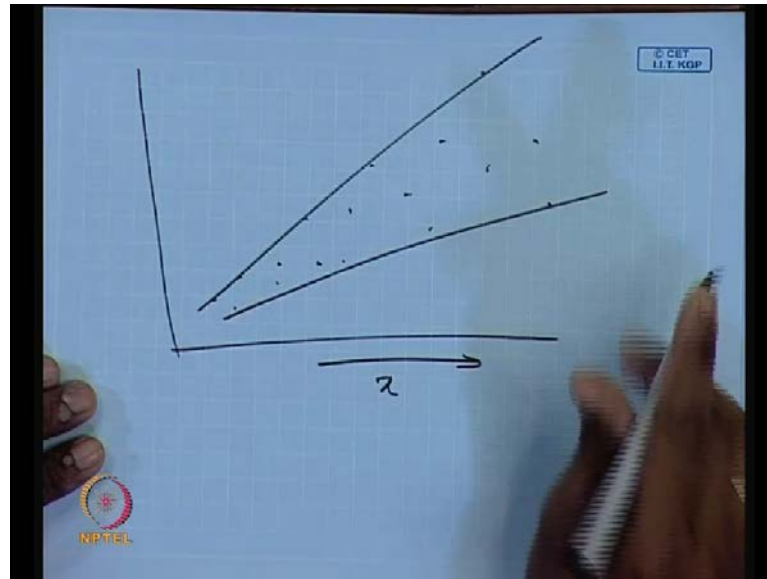
- When conditional variance about the regression line is a function of the independent variable, it may be expressed as :
$$\text{Var}(Y|x) = \sigma^2 g^2(x)$$
where $g(x)$ is a predetermined function, and σ is an unknown constant.
- Assumption : Data points in regions of small variance have more 'weight' than those in regions large variance. So, we assign weights (w_i) inversely proportional to the variance.

 Probability Methods in Civil Engineering: H3L4 Dr. Rajib Maity, IIT Kharagpur 21

So, now we will deal with the regression with the non constant variance now when the conditional variance about the regression line is a function of independent variable, it may be expressed as variance of Y given x is equal to sigma square multiplied by g square x . Now, this $g x$ basically is the predetermined functions. Some function is that, how it is varying and this should be multiplied with the sigma square and when it is variance you know that any function or constant that is multiplied with this variance so that, we make it square.

That is, we discussed in the earlier lectures. So, this is that sigma square g square x . So, now, this sigma is an unknown constant and here the assumption is that data points in the region of this small variance have more weight and than those in the region of this large variance. So, we assign the weight w_i inversely proportional to the variance. So, some weight we have to put and our assumption is that when the data is having the small variance.

(Refer Slide Time: 36:33)



Now, if you see this diagram basically, if I say that this is varying means, suppose this what we can see in this literally. So, we can easily see that as it is going and so, if this the x as this x is varying basically the range is changing. So, here the weight will be more in this zone where the variance is less and here the weight will be less where the variance is more basically that is what. So, in this way it is inversely proportional to the variance.


(Refer Slide Time: 37:03)

Regression with non-constant variance

$$\hat{w}_i = \frac{1}{\text{Var}(Y | x_i)} = \frac{1}{\sigma^2 g^2(x_i)}$$

■ The squared error is calculated as :

$$\Delta^2 = \sum_1^n \hat{w}_i (y_i - y_i')^2 = \sum_1^n \hat{w}_i (y_i - \alpha - \beta_{\frac{1}{g}} x_i)^2$$



Probability Methods in
Engineering: EE314

Dr. Rajib Maity, IIT Kharagpur

22

So, this how the weights are given. So, $1/\text{variance of } Y \text{ given } x_i$ which is the $1/\sigma^2 g^2(x_i)$ the squared error is calculated as this sigma square is equals to that this

weight is weights we will put and then that your that difference square and sum it up. So, this is the equals to 1 to n.

(Refer Slide Time: 37:33)

Regression with non-constant variance

- To find the least-squares estimates of α and β , the total error is minimized and thus we obtain α and β , as :

$$\hat{\alpha} = \frac{\sum_{i=1}^n w_i y_i - \hat{\beta} \sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

NPTEL Probability Methods in Engineering: H314 Dr. Rajib Maity, IIT Kharagpur 23


So, now this y_i cap again that estimate will get from this alpha minus beta x_i now to find that least square estimate of alpha and beta the total error is minimize and thus we obtain this alpha and beta and following the same principal that we have discussed for this constant variance we will just get that error and error is partial derivative is taken equated to 0 and after solving those equation we will get the estimate of alpha is equals to through this expressionsay w_i into y_i minus beta cap of this $w_i x_i$ divided by w_i summation of all this w_i .

(Refer Slide Time: 38:08)

Regression with non-constant variance

$$\hat{\beta} = \frac{\sum_{i=1}^n w_i (w_i y_i x_i) - \left(\sum_{i=1}^n w_i y_i \right) \left(\sum_{i=1}^n w_i x_i \right)}{\sum_{i=1}^n w_i \left(\sum_{i=1}^n w_i x_i^2 \right) - \left(\sum_{i=1}^n w_i x_i \right)^2}$$

where $w_i = \sigma^2 w_i' = \frac{1}{g^2(x_i)}$



Probability Methods in
Engineering: EE715

Dr. Rajib Maity, IIT Kharagpur

24

And, beta cap will be obtain through this expression even though this expression looksthroughlike a little bit cumbersome, but,thing is that this is we get following the same principal that we havedone for this constantvariance only thing here the one that weight function is coming and which the weight you can see that this weight is equal to sigma square wi prime and this sigma square.Ifwe just multiply whatever the equation that we have used here.


(Refer Slide Time: 38:36)

Regression with non-constant variance

$$\hat{w}_i = \frac{1}{Var(Y | x_i)} = \frac{1}{\sigma^2 g^2(x_i)}$$

■ The squared error is calculated as :

$$\Delta^2 = \sum_1^n \hat{w}_i (y_i - y_i')^2 = \sum_1^n \hat{w}_i (y_i - \alpha - \beta x_i)^2$$



Probability Methods in
Engineering: EE715

Dr. Rajib Maity, IIT Kharagpur

25

So, this σ^2 will be cancelled. So, it will be $1/g^2$. So, now this g^2 generally same function we will use and that function of this x should be there to when we are determining this w_i to get the estimate of this α and β and the conditional variance is calculated as $S^2_{Y|X}$ is the s^2/g^2 and S_{xi} – this is the standard deviation, the square root of this positive square root of this. So, s multiplied by the g ; you can see here that this conditional standard deviation is a function of that x where this S here is that summation of $(y_i - \alpha - \beta x_i)^2$ divided by $n - 2$.

(Refer Slide Time: 39:30)

Example

Q. The maximum settlements and maximum differential settlements of a 10 storage tanks is as shown in table. The differential settlement appears to increase with maximum settlement. Assume that the conditional standard deviation of the differential settlement Y increases linearly with the maximum settlement X , or $\text{Var}(Y|x) = \sigma^2 x^2$. Obtain the regression equation for estimating the expected maximum differential settlement Y on the basis of information for the maximum settlement X of the tank.

NPTEL Probability Methods in Civil Engineering Dr. Rajib Maity, IIT Kharagpur 26

Now we will take one example on this one. It will be more clear in that way where the variance is dependent on the value of X . The maximum settlement and the maximum differential settlement of 10 storage tanks, this is wrong; of 10 storage tanks is as shown in the table. The differential settlement appears to increase with the maximum settlement assumed that the conditional standard deviation of the differential settlement Y increases linearly with the maximum settlement X or this is what is told. That is, linearly it increases; that means, that g x that the function that we have told this is g x is equals to X .

So, the variance of Y given X is equals to $\sigma^2 x^2$; that function square obtained, the regression equation for estimating the expected maximum differential

settlement y on the basis of the information for the maximum settlement of X of that tank to do this one; this is the data that 10 different data set is given here.

(Refer Slide Time: 40:39)

Example...Contd.

Tank No:	Maximum Settlement (cm) x_i	Maximum Differential Settlement (cm) y_i
1	0.32	0.3
2	0.5	0.8
3	0.8	1.1
4	0.9	0.6
5	0.8	1.0
6	1.2	1.3
7	1.3	1.5
8	1.1	1.1
9	1.5	0.7
10	1.6	0.8

NPTEL Probability Methods in Civil Engineering: Dr. Rajib Maity, IIT Kharagpur 27

So, this is the maximum settlement there is a maximum differential settlement. So, here we have to regress that maximum differential settlement on this maximum settlement. So, our variable here in this following the notation that we have used is the this is of our y and this is our x.

(Refer Slide Time: 41:05)

Example...Contd.

Sol.:

The conditional standard deviation of the differential settlement Y increases linearly with the maximum settlement X, or

$$\text{Var}(Y|x) = \sigma^2 x^2$$

So

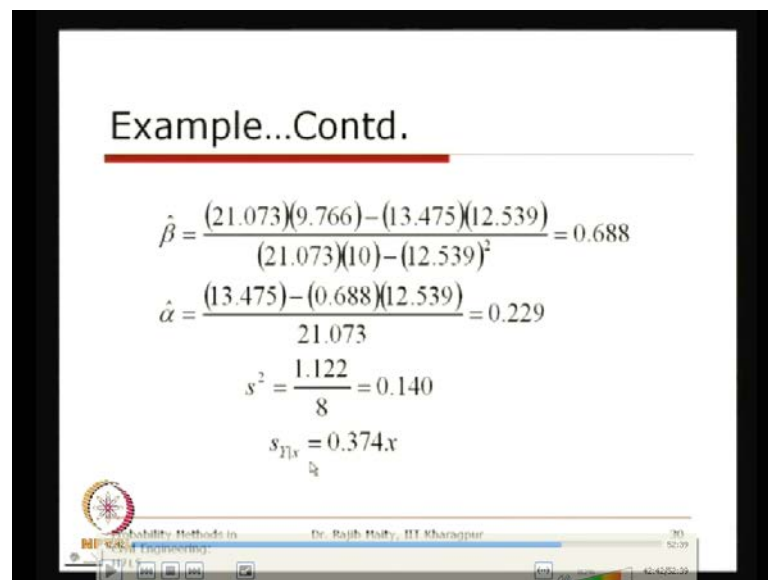
$$w_i = \frac{1}{x_i^2}$$

NPTEL Probability Methods in Civil Engineering: Dr. Rajib Maity, IIT Kharagpur 28 41:35(41:39)

So, the conditional standard deviation of the differential settlement Y increases linearly with the maximum settlement X or variance Y on condition x is equal to $\sigma^2 x$'s square.

So, this is the relationship that is given. So, this x this function actually is predetermined as we have seen in that theory. So, here that w_i is the inverse of that function. So, $1/y_i$ square. So, this is the weight age that is the w_i and with that w_i , for all these we will get this weight and basically, this is input. This is also we know the observed data this is the weight which is an inverse to this x_i and. So, these things we can calculate $w_i x_i$, $w_i y_i$, $w_i x_i^2$ and $w_i y_i^2$.

(Refer Slide Time: 42:09)



Example...Contd.

$$\hat{\beta} = \frac{(21.073)(9.766) - (13.475)(12.539)}{(21.073)(10) - (12.539)^2} = 0.688$$

$$\hat{\alpha} = \frac{(13.475) - (0.688)(12.539)}{21.073} = 0.229$$

$$s^2 = \frac{1.122}{8} = 0.140$$

$$s_{y|x} = 0.374x$$

The slide also shows a presentation footer with the text 'Probability Methods in Engineering', 'Dr. Rajib Mallik, IIT Kharagpur', and a slide number '20'.

So, if we use this one and then up to this of this table we can calculate and based on this we will estimate that alpha and beta and here the beta is estimate of this beta is 0.688 and estimate of this alpha is 0.229 and the S^2 is 0.140 and this standard deviation of y given x equals to $0.374x$.

(Refer Slide Time: 42:45)


Example...Contd.

- The regression equation for the estimating the expected maximum differential settlement Y on the basis of information for the maximum settlement X of the tank is as:

$$E(Y|X) = 0.229 + 0.688x$$

- Standard deviation is

$$s_{Y|X} = 0.374x$$

 Probability Methods in
Civil Engineering: II/15

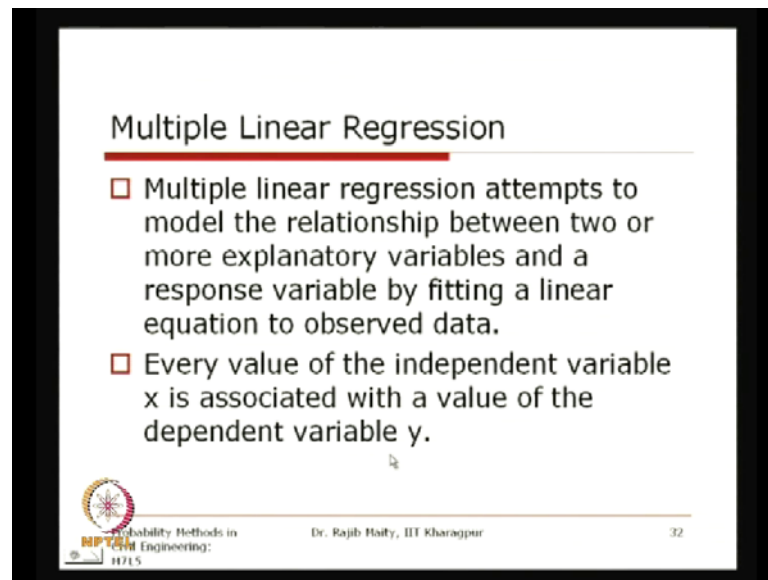
Dr. Rajib Halty, IIT Kharagpur

31

So, you can see that as x increases this standard deviation also increases which is the function of this x of the variable x . Now, the expected value of this y is the regression equation for the estimating the expected maximum differential settlement Y on the basis of information for the maximum settlement X of the tank is as expected value of this y given x is equals to 0.229 plus $0.688 x$.


So, this is that expected value of y given x and this is expected value of this is the standard deviation of y given x and now, using this relationship basically when how we are getting this 0.345 here that we have seen that S_{square} is this one. Basically, we are using this α and β estimate to calculate this one first and this total we are getting this is the sum of square error weightage sum of square error and from there we are getting this s_{square} and from there it we getting that given that variance, sorry, standard deviation of y given x $0.374 x$.

(Refer Slide Time: 43:44)



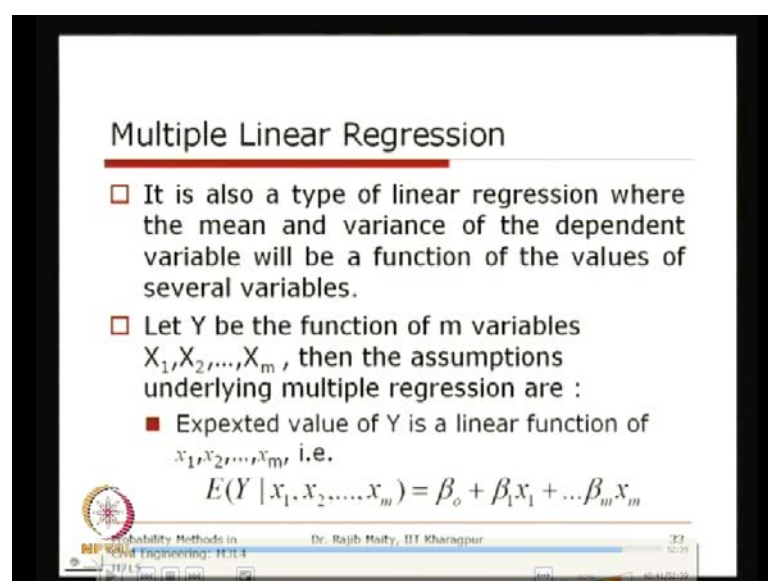
Multiple Linear Regression

- Multiple linear regression attempts to model the relationship between two or more explanatory variables and a response variable by fitting a linear equation to observed data.
- Every value of the independent variable x is associated with a value of the dependent variable y .

 Probability Methods in Engineering: H215 Dr. Rajib Maity, IIT Kharagpur 32


So, next we will take that multiple linear regression and here you know that. So, far whatever we have discussed it is that regression and one dependent variable, one response variable one target variable was there. Now, in case when we are having that more than one random variable then we have to go for this more than one dependent variable then we have to go to this multiple linear regression. So, this multiple linear regression attempts to model the relationship between two or more explanatory variables and a response variable by fitting a linear equation to the observed data. Every value of independent variable is associated with a value of the dependent variable y .

(Refer Slide Time: 44:42)



Multiple Linear Regression

- It is also a type of linear regression where the mean and variance of the dependent variable will be a function of the values of several variables.
- 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.
$$E(Y | x_1, x_2, \dots, x_m) = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m$$

 Probability Methods in Engineering: H214 Dr. Rajib Maity, IIT Kharagpur 33

It is also a type of linear regression where the mean and variance of the dependent variable will be a function of values of these several variables. So, here instead of that using that one independent variable that is the X ; that X and our dependent variable was a Y earlier case.

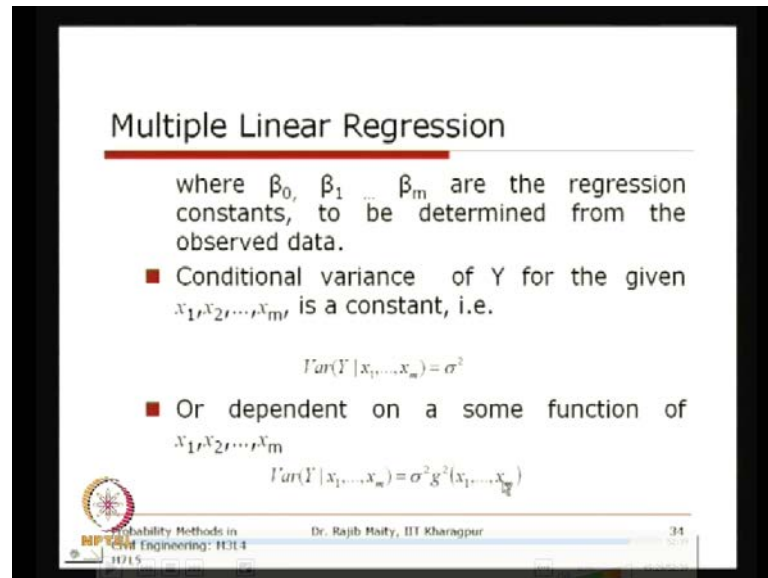
So, here what you can see is that Y be the function of m variables instead of only one. So, far what we have discussed here is, Y is the function of m variables which is x_1, x_2 up to x_m then the assumptions underlying the multiple regression are the expected value of Y is a linear function of x_1, x_2 up to x_m that is the expected value of y given the information of this independent variable x_1, x_2 up to x_m equals to that $\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m$. One thing I will just, one correction I will just do before I proceed further in this linear regression - with one between x and y . So, where only one input was there I might have sometime mentioned that this x is in this expression when we are regressing y on x , I might have sometime mentioned that this x is your dependent variable and y is the target variable.

So, the correction will be that x is your independent variable and y is your dependent variable. Sometimes, for the y we can mention that this is the target variable, response variable, dependent variable and all and basically, when we are referring to this x this is the independent variable. Earlier in this case, when we were discussing the simple regression that time, only one dependent variable was there. Now, what we are discussing here is that, we are having more than one dependent, more than one independent variable to model that dependent variable y and this is through a linear function which is that $\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m$. Now, basically how the concept is taken through is that.

Now, for that when you see that there is only one independent variable and one dependent variable between x and y , we are basically fitting a straight line through the observed data point. Now, when we are having more than one input say for example, if I just say there are 2 inputs x_1 and x_2 and our target - our dependent variable is y then basically, you can visualize, you can conceptualize in this way that this is a 3 dimensional space over which the two axis is one for the x_1 , other for the x_2 and we are basically fitting one surface; one straight, one linear surface through the data point in the 3 dimensional space. So, this is in case of when there are two independent variables and one dependent variable. Now, similarly you can extend it to the higher dimension and this; so that, for

the independent variable, the relationship is generally that $b_0 + b_1x_1 + \dots + b_mx_m$. Now, we will follow the similar procedure to estimate these parameters as well. That is, we have to first find out what is the error and that error should be squared up; sum them. So, there is the sum of square error then it is minimized with respect to the parameters to get those expressions.

(Refer Slide Time: 48:44)



Multiple Linear Regression

where $\beta_0, \beta_1, \dots, \beta_m$ are the regression constants, to be determined from the observed data.

- Conditional variance of Y for the given x_1, x_2, \dots, x_m is a constant, i.e.

$$\text{Var}(Y | x_1, \dots, x_m) = \sigma^2$$
- Or dependent on a some function of x_1, x_2, \dots, x_m

$$\text{Var}(Y | x_1, \dots, x_m) = \sigma^2 g^2(x_1, \dots, x_m)$$

Probability Methods in Civil Engineering: H3E4
Dr. Rajib Maity, IIT Kharagpur
34

So, where this $\beta_0, \beta_1, \beta_m$ are the regression constant to be determined from the observed data and the conditional variance of Y for the given x_1, x_2 up to x_m is a constant that is variance of this Y given this input is equal to sigma square. Or, this is in case of when it is constant or it may be dependent on some function of this x_1, x_2, x_m . So, when it is dependent, when it is varying when it is non constant as we have used in the simple regression case. So, this variance of y given x_1, x_2, x_m is equals to sigma square multiplied by the square of some function of this x_1, x_2 up to x_m .

(Refer Slide Time: 49:31)

Multiple Linear Regression

- The regression analysis determines the estimate for $\beta_0, \beta_1 \dots \beta_m$ and σ^2 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 Engineering Dr. Rajib Maiti, IIT Kharagpur 35

Now that expression that is the regression analysis determines the estimates for this beta naught, beta 1, beta m and the sigma square, sorry, this will be square and the sigma square based on the given data $x_{11}, x_{12}, x_{13}, \dots, x_{1n}$ up to x_{mi} and i is varying from the 1 to n . So, we are having then n set of i can say that n set of data that is $y_1, x_{11}, x_{12}, x_{13}, \dots, x_{1n}$. Similarly, i will have another set of this data. So, there are, m is the number for the number of the dependent variable and n is the number of what? How many sets of the observed data that is available to us.

So, based on this we can, whatever the expression that we have seen in this expression can be slightly modified as this one. That is, the α plus $\beta_1(x_{11} - \bar{x}_1)$ plus up to this that $\beta_m(x_m - \bar{x}_m)$. So, how we get this one is that, this \bar{x}_1 is the mean of whatever we have seen in this x , in the variable of x_1 and the \bar{x}_m is the mean of that observed data of the, for the independent variable x_m . So, basically what we are replacing is that this constant beta naught is basically a adjusted is basically replaced. So, this alpha you can see that this alpha is equals to that beta naught plus $\beta_1 \bar{x}_1$ plus $\beta_2 \bar{x}_2$ plus up to this $\beta_m \bar{x}_m$.

So, here what we can get is that, from this expression we can estimate that is alpha, beta 1, beta 2, beta m and from that estimate of this alpha and obviously, beta 1, beta 2 of the same. If we put it here, we will get what is the estimate for this beta naught and this one we will see. We will continue from this point onwards in our next lecture and what is in this

linear regression part, what we have seen in today's lecture is the linear regression and linear regression with respect to the constant variance and some or the non constant variance. We have seen one example for each case and this. So, in the next class what we will see is that multiple linear regression and then we will also see the non-linear regression and we will see the R^2 relation as a measure of how strong the relationship is captured; that we will see in the next class; thank you.