**Data Mining**
**Prof. Pabitra Mitra**
**Department of Computer Science & Engineering**
**Indian Institute of Technology, Kharagpur**

**Lecture - 38**
**Regression II**

We continue the discussion on linear regression.

(Refer Slide Time: 00:23)
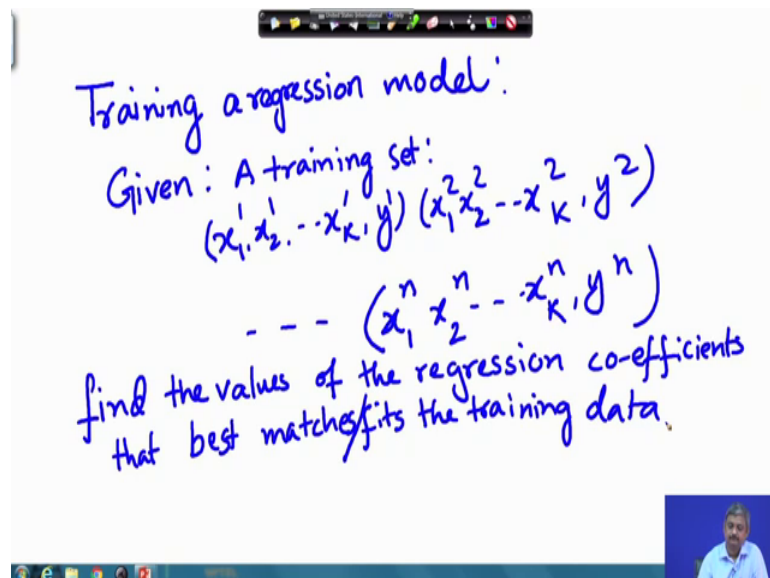


(Refer Slide Time: 00:25)

As we discussed earlier we have a regression model of the following form. In general we can write as and then independent valuables.
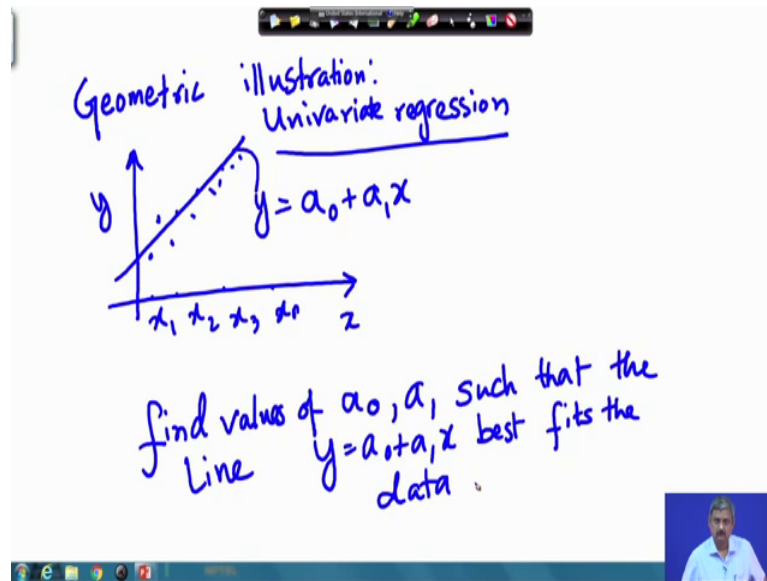
(Refer Slide Time: 02:38)



As (Refer Time: 03:51) so, this is the illustration.
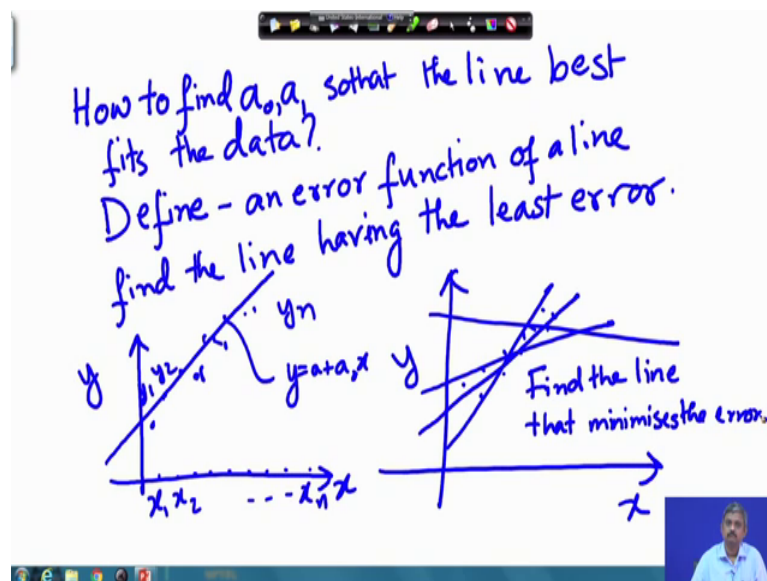
(Refer Slide Time: 03:51)
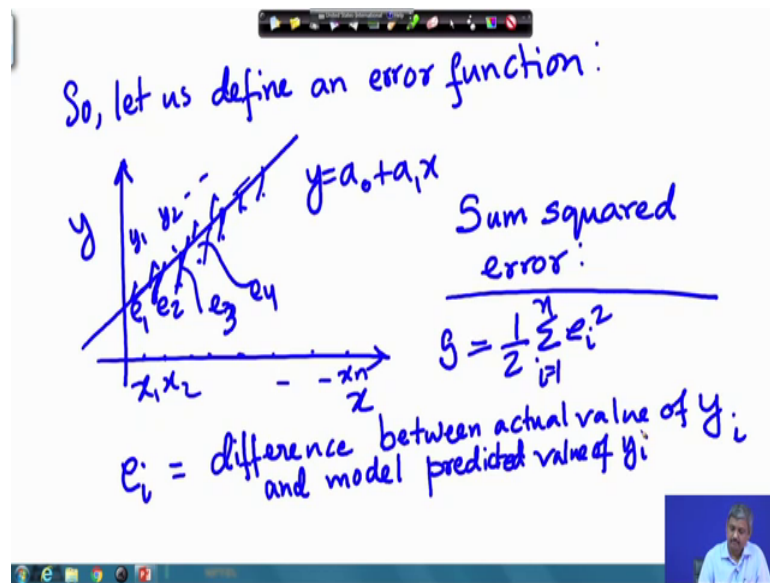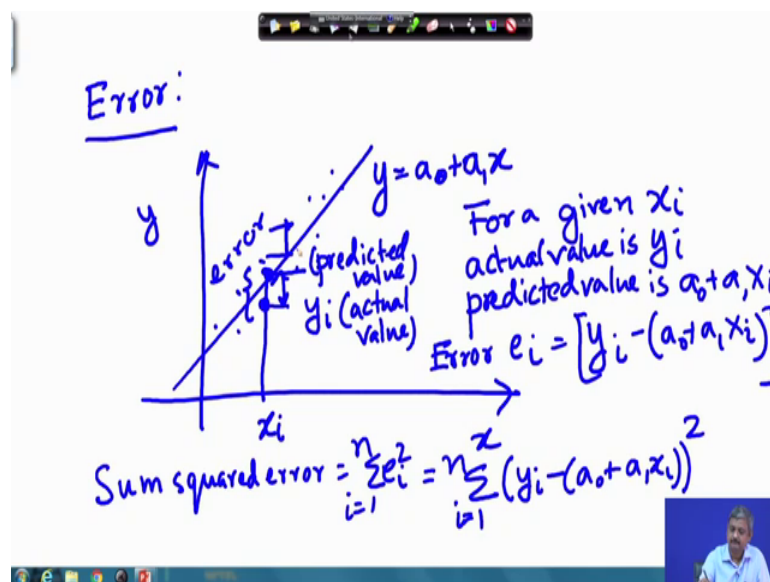
Data, find a 0 a 1.

So, that it best fits the data. Ok define some error function and minimise it. There will be some error it will not exactly match the points. So, try to find the line.

(Refer Slide Time: 09:53)



This line, this line or this line or this line; So, what is the error? You measure how much this point is off from the line. Ok square up all these add them up. Ok, square these are the error sum them; the difference between actual value and predicted value.

(Refer Slide Time: 12:44)



So, I xi this is the error, this defines why I am taking square because see error can be positive as well as negative moment I take both are equally bad.

(Refer Slide Time: 15:38)



So, moment I make it positive a square it is becoming positive, such that minimised ok. So, find a 0 a 1 a 2 so that this quantity is smallest for a given training set ok; So, here I have used a slightly different methodology this is alpha and these coefficients are called the model parameters; So, the squared.

(Refer Slide Time: 18:03)

What S? It is a function of the parameters what if you choose certain parameters sum error you will get. So, this is my error, where like this. So, I can write this like this e transpose into e is this thing there is this thing. So, if we expand you get this.

I am just writing it down so finally, you will get theta is the spectre. Now, at minima this will be equal to 0 the minima will be derivative; minima this thing ok.

(Refer Slide Time: 21:42)



(Refer Slide Time: 22:40)



So, there are different methods, get theta there are different methods ok.

(Refer Slide Time: 23:49)



So, this is like solving simultaneous equations ok. So, you can extend it to multivariate case also.

(Refer Slide Time: 23:57)



Alright so, that is how ok.
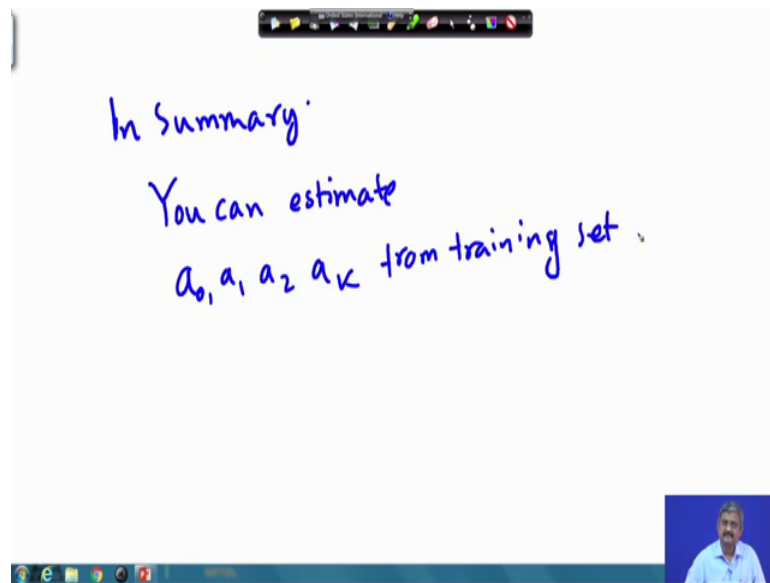
So, in the next lecture I stop here, in the next lecture I will explain how to extend it to non-linear cases.

Thank you.