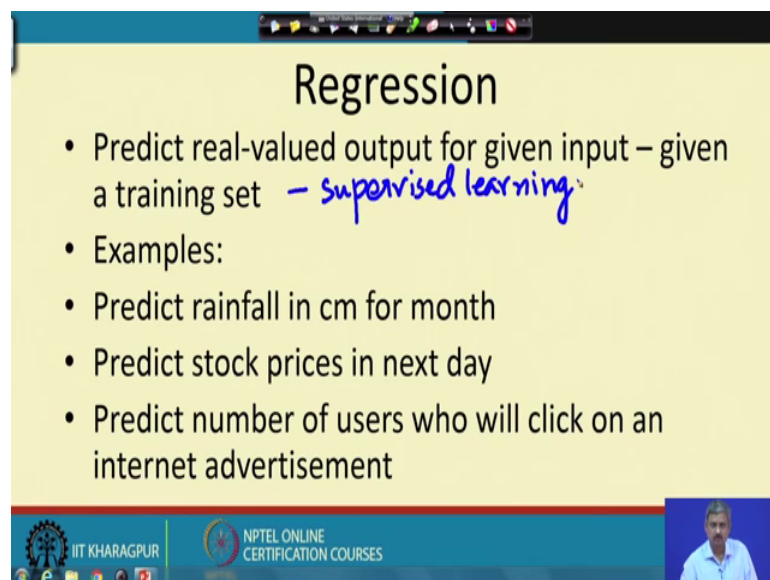**Data Mining**
**Prof. Pabitra Mitra**
**Department of Computer Science & Engineering**
**Indian Institute of Technology, Kharagpur**

**Lecture - 37**
**Regression – I**

We will now discuss another Data Mining problem known as Regression Analysis. So, in this week, the 8th week, we will discuss Regression Analysis along with few other techniques that are essential for Data Mining.

(Refer Slide Time: 00:40)



So, we had previously studied mainly two important Data Mining class, besides the association rule mining. One is the classification, where there are some predefined set of classes. For example, a spam or a non-spam; for example, a news article falling in a category say sports, politics, local. So, there are defined categories and you have a training example, a set of training example, where there is a supervised output; that means, for a every example, what should be the category, it is mentioned.

(Refer Slide Time: 01:45)



So, if I quickly write down, a set of predefined categories or; you have a training set where not only the attributes of the examples but also the class levels are available you have to predict the class level for a new example. So, you have examples of spam and non-spam email, you have to predict a new email whether it is spam or not a spam. So there are two aspects, one is there is this class level information along with 10 attributes which is called Supervised learning. Second is, you are predicting for an unknown example, why it is predictive mining ok. So, these are the two important problem.
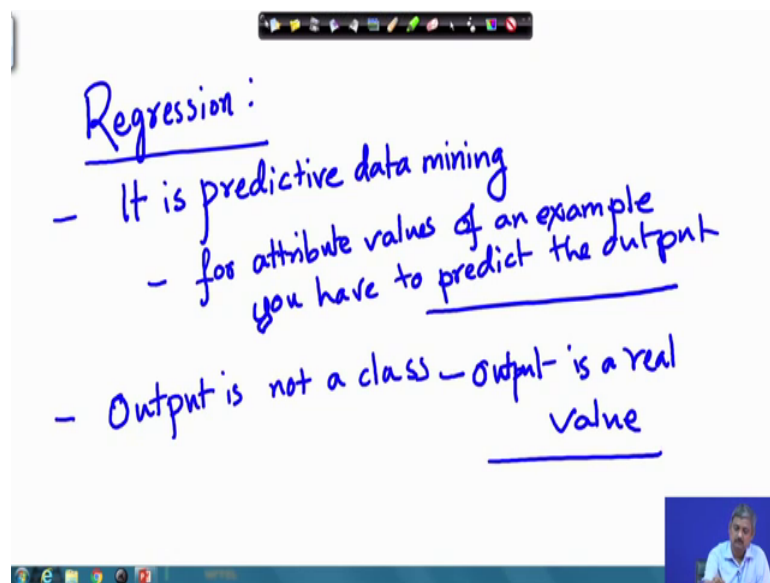
(Refer Slide Time: 04:27)

Similarly, we also saw another task. We also saw the Clustering task where no predefined classes. You just find to try to find the homogenous groups in that data. So, it is that is why called exploratory data mining and training examples, they have attribute values only; no class information is available ok. So, nobody tells that for this example, this should be the class. So, there is no supervision. So, it is unsupervised learning. That is exploratory and unsupervised.

(Refer Slide Time: 07:02)



Now, we will look at another problem, Regression. It is predictive data mining; that means values of an example. You predict output but output is not a class; it is not a class. This, if it is a class, we can probably represent them by replace them by in tiers. Say spam is 0; non-spam is 1 ok, fraud is 0; non-fraud is 1, fraud is minus 1, non fraud is we can represent them by a, but here it is not a class. So, output is no longer an integer, output is a real value ok.

So, regression is predictive mining where output is real and you still have a training set. So, it is supervised also, learning also ok. So, some examples, may be you have gone to weather dot com and you want to predict the rainfall let us say in centimetre for the next one month, some real value say 230 centimetre or something. Or for example, you want to predict the value of a stock. In the next step that is also a real value whereas, if you predict if the stock goes up or goes down that is a 2 class classification problem ok.

Another example, many people study is the number of users who will or the fraction of users who will click on certain internet advertisement.

So these are all examples of Regression.

(Refer Slide Time: 10:30)



One of the most common technique of Regression is what is called a Linear Equation ok.

(Refer Slide Time: 10:55)



So, let me explain what it is. So, what is given to me is, at a set of values of x, x 1, x 2, x 3, x 4 up to say x n, I have the corresponding output y 1, y 2. So, x i say for example, say

is the temperature today of certain place and y i can be volume tomorrow, ok; rainfall volume tomorrow all right. So, there can be other also. Let me take another example. So, x i is let us say, say how does traffic density, number of vehicles in a road depend on temperature?
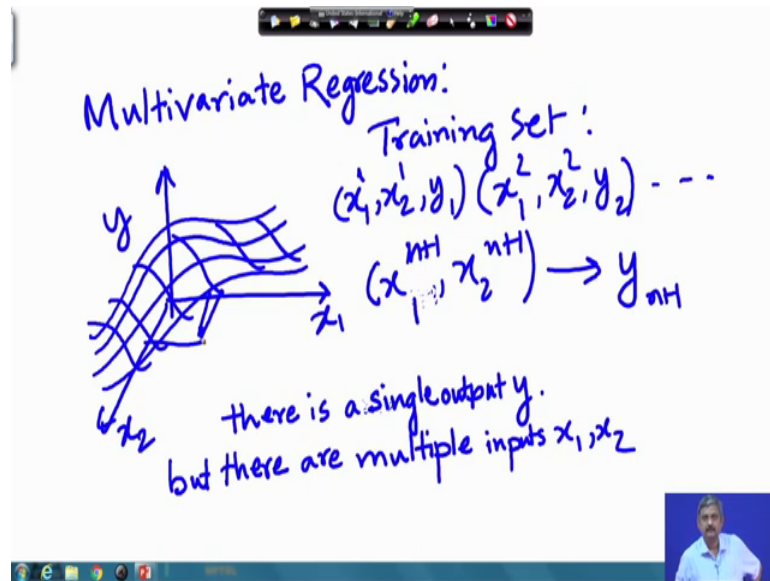
If it is very hot people do not go out, if this very cold people do not go out. So, can I predict the traffic density as a function of the temperature of the day and that is a Regression problem ok.
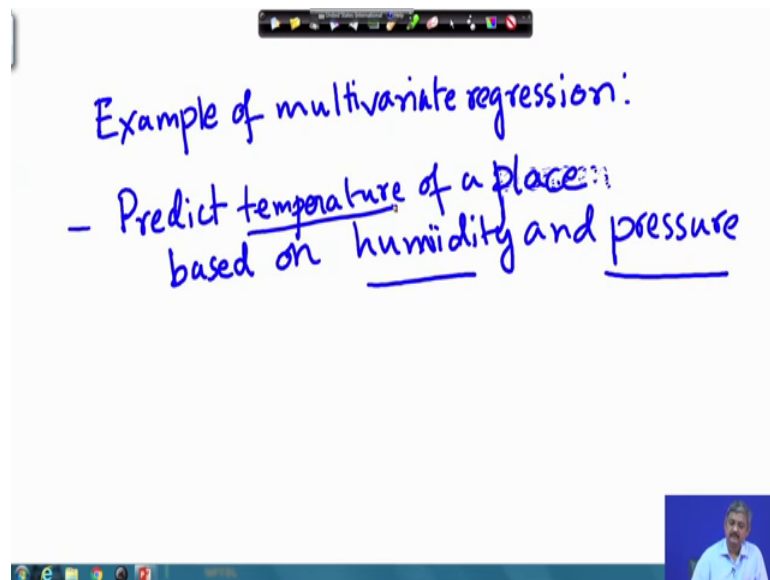
(Refer Slide Time: 14:57)



So, in general, Training Set consists of y n. Given a new point x n plus 1, what should be the value of y n plus 1? Ok, that you have to predict, all right. So, may be you extend the curve like these and see what value it takes, like that you do ok. So, this particular example when the x has only a single value x 1, x 2 is called a Univariate regression.

(Refer Slide Time: 17:10)



We can have, we can have Multi value. There is a x 2 ok. So, here it what will be and so, on. And for a new point sorry 1, you have to predict is Multi value.

(Refer Slide Time: 19:44)



So, basically one point here, you know you have to predict it. For 2, given x 1, x 2, you have to predict 1 ok. So, let us see some example of Multivariate regression. Place based on pressure, based on humidity and pressure, you predict the temperature place. So, there are two inputs, bivariate and one output. Actually, in general, there can be multiple outputs also.

(Refer Slide Time: 21:21)



So let us write down what is called a Regression Model. Regression Model is nothing but expressing y as a function of the input x ok. As a function, if it is univariate. So, these are the attributes and this is attribute, for multivariate case multiple attribute, univariate one attribute; y is called the output or the dependent variable; whereas, input or the independent variables and f is called the function or the regression model ok.

(Refer Slide Time: 23:42)



So, the idea is that dependent variables or idea is, f determines how ok. That is the idea.

So you basically, f is this function, f is this function.

So, what is Linear Regression? When f is a linear function ok; so, for linear regression, some constant alpha 1, alpha 2 or not alpha; let me write the mouse all right, write them as a 1, a 2. So, f is this linear function, f is this linear function are the regression are the called the regression Co-efficient; a 1, a 2, a k, a naught, they are called the regression Co-efficient ok.
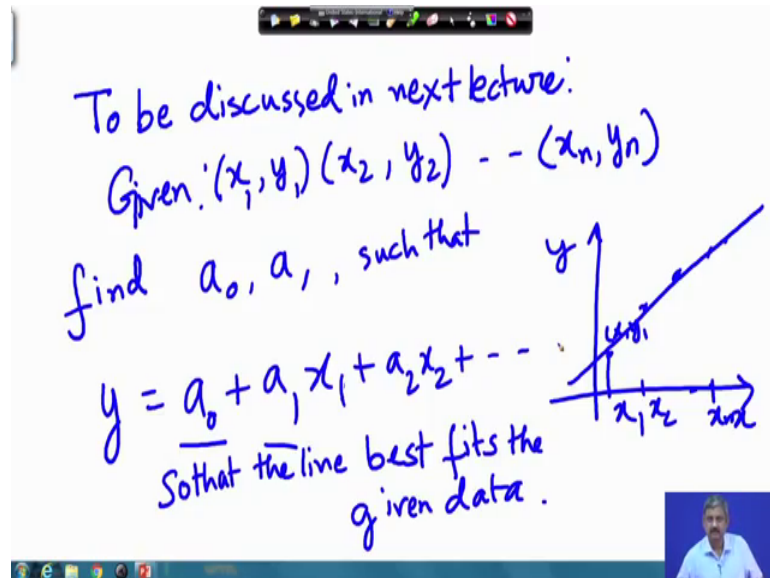
So, just f is of this particular form, a 1, a 2 a simple ok. So, in the univariate case, defining values if they are given, then a 1 x, straight line; a naught is the shift and a 1 is the slope. In bivariate, it is a plane.

(Refer Slide Time: 29:15)



So that is linear equation as you can expect. If we are given some regression points say x 1, y 1; x 2, y 2; given some points, a naught such that so, basic. So, you find values of a naught and a 1, so that the given data it matches; best fits the given data that we will find out. If we have higher dimension, we have more terms. So, the all these terms you have to find out ok.

So, that is the Regression design problem; so that we will do in our next lecture, how to find that.

Thank you.