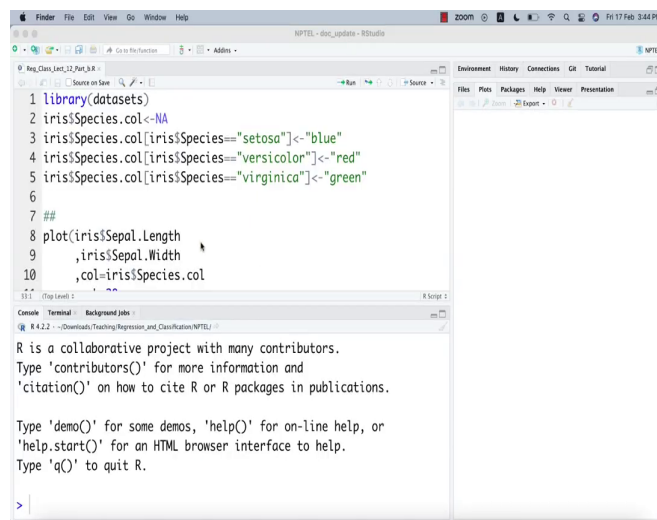**Predictive Analytics - Regression and Classification**
**Prof. Sourish Das**
**Department of Mathematics**
**Chennai Mathematical Institute**

**Lecture - 41**
**Hands on with R with Iris dataset**

Welcome to the part b of lecture 12. In this video, we are going to do some Hands on. So, in this hands on we are going to load the dataset library which has the; which has the iris dataset.
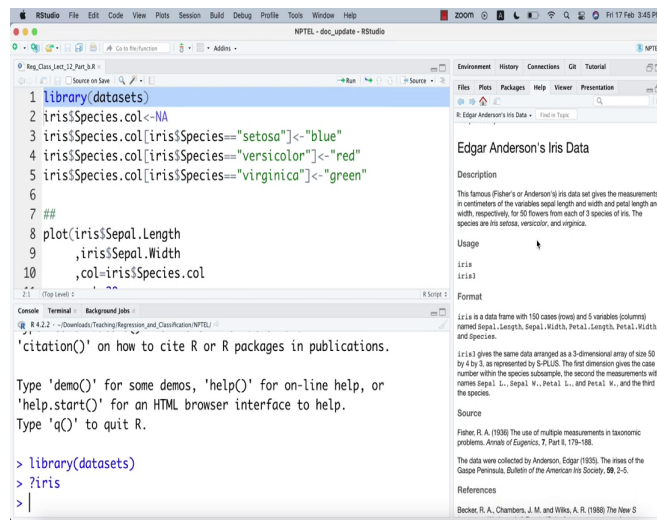
(Refer Slide Time: 00:23)

(Refer Slide Time: 00:32)



So, if you go to the load the, you write iris. So, Edgar Anderson's iris datasets, it is also known as Fisher's iris dataset. So, it has iris is a English flower, which has three subspecies, one is called Setosa, one is called versicolor, and one is called virginica.

(Refer Slide Time: 01:02)



In this dataset, what happens is if you just say, let us say, iris, head iris. So, it has four predictors, sepal length, sepal width, petal length, and petal width. This four based on these four phenotype can you say which species the flower belongs to? Ok. So, what we I have done here, I for different colors setosa, versicolor, and virginica.

I have given different color, for different species, I have now had different color. So, now, if you go to head iris.

(Refer Slide Time: 01:55)



So, you will see that in along with that, let me just.

(Refer Slide Time: 01:59)



Yeah. So, now, this is the new column, that last column is a new column. For setosa, I had blue, versicolor I have given red, and virginica I given green.

(Refer Slide Time: 02:18)



Now, if I now if you run this piece of code, I just plotting from iris dataset, I am just taking the sepal length and sepal width. I am just extracting it and giving it as a x value and y values. And if you run this plot, then you got this plot, you want, you can give a grid also.
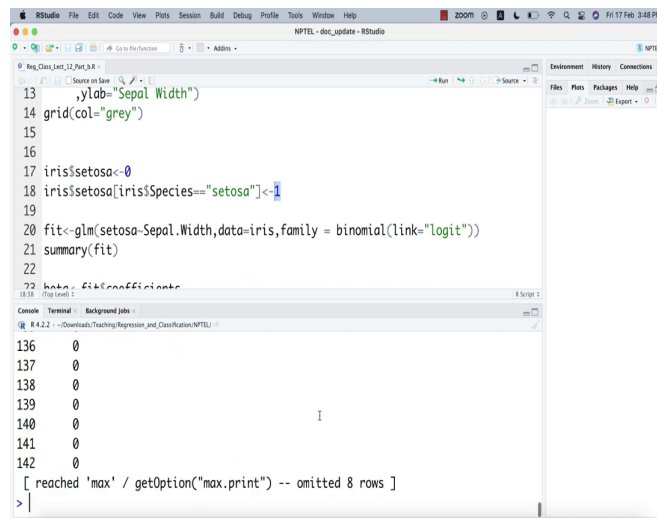
So, you can see these are like bluer setosa, we have given bluer setosa, reds are red points are all versicolor, and the green points are virginica.

(Refer Slide Time: 02:58)



Now, what I have done in that dataset, I have created a one-hot encoding for setosa. If it is setosa, then it will have a 0 or it is 1.

(Refer Slide Time: 03:21)



So, for all it is 0, but whenever it will find setosa, it will get a 1 otherwise, it will be 0. So, now the dataset, let me just let me just show you how the dataset look likes now, iris, ok.
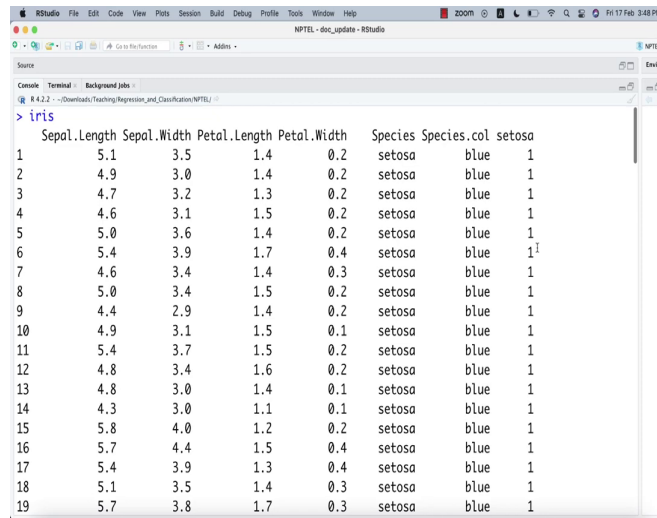
(Refer Slide Time: 03:36)

(Refer Slide Time: 03:40)

(Refer Slide Time: 03:42)



So, let me just, yeah. So, first few values are setosa and then when it is versicolor, it is coded as 0, when it is virginica, this is coded as 0, ok.

(Refer Slide Time: 04:25)



And the column name is setosa. So, that is how I created a one-hot encoding or binary class variable, sometime it is called indicator variable. So, let me run the model, glm, setosa as a function of Sepal Width, data equal to iris, you give family equal to binomial link equal to logit.

(Refer Slide Time: 04:29)

(Refer Slide Time: 04:32)



And if you run summary fit, then this is the fit that you will get.

(Refer Slide Time: 04:35)



If you run the beta, you can from the fit, you can extract the coefficients. Now, in the beta, I have the coefficient. Now, what I am going to do, I am going to calculate eta X matrix with 100 values with which X 2 takes values with minimum value of Sepal Width to maximum value of Sepal Width.

(Refer Slide Time: 05:04)



So, if you just run it. So, minimum is I think 2.2 is the minimum.

(Refer Slide Time: 05:08)



And max is the 4.4. And in between, it just fill up with some values in a equal width. And then I calculate the eta, these are the eta values or z values, latent variable values. And then I calculate the probability.

And then I plot the X 2 versus Sepal Width versus sorry, but that sepal width versus the probability of setosa. So, if the Sepal Width increases, clearly probability that the flower is setosa increases. So, you can put a grid also in this. So, that is how we, this is that is how we can, you can draw probability of p, you can plot p against some predictor values.

So, I will stop here, but this iris data set shows you that your target variable could be not necessarily has to be binary class. It could be multi class variable. And then in that case, you have to you, you have a multi class classification. Because you have three class here. Remember that you have to do three class, not binary class.

You have a setosa, versicolor and virginica, three sub species are there. When you have more than two class, basically it is a k class problem or multi class problem. The most popular and the old and tested method is linear discriminant analysis.

So, in the next video, we are going to start linear discriminant analysis or Fisher's linear discriminant analysis. So, for now, thank you very much for your attention. See you in the next video.