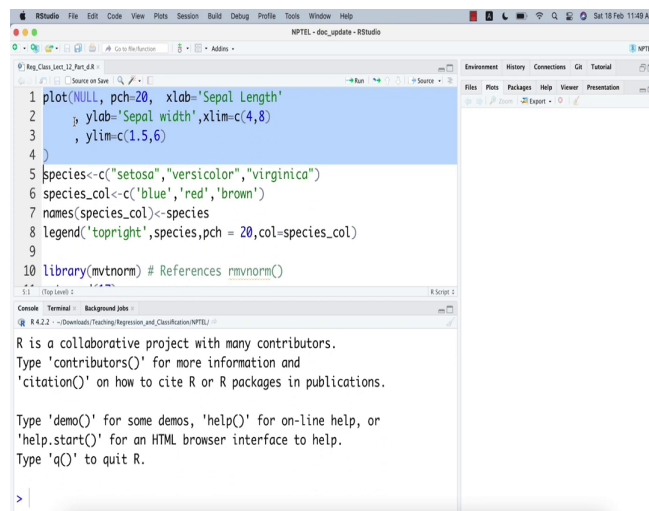**Predictive Analytics - Regression and Classification**
**Prof. Sourish Das**
**Department of Mathematics**
**Chennai Mathematical Institute**

**Lecture - 43**
**Hands on with R: Implement LDA**

Welcome back to the part d of lecture 12. In this video, we are going to do some Hands on.
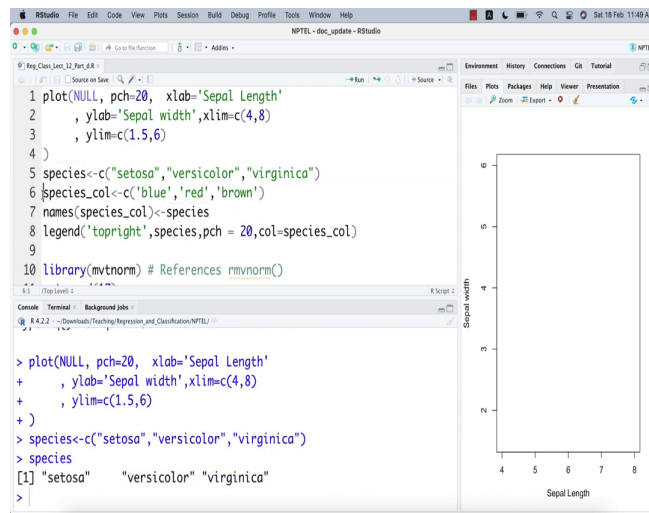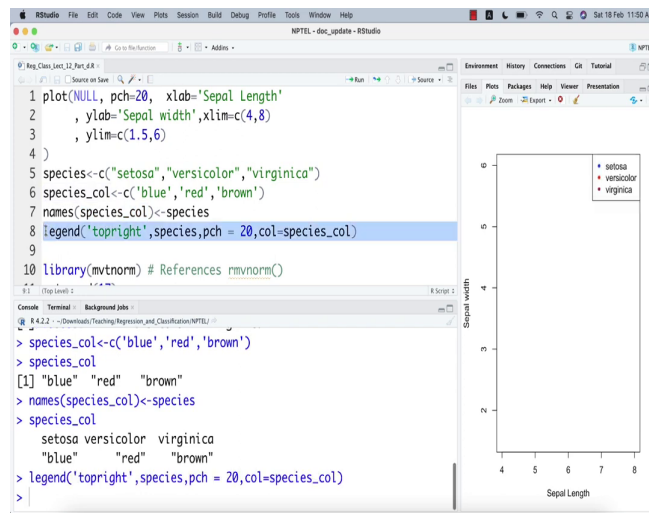
(Refer Slide Time: 00:34)



So, let me start my R; in this. So, we will still we will work with the iris dataset. So, here I am going to first create a null plot. So, I we will just first see how the you know how the data looks like.
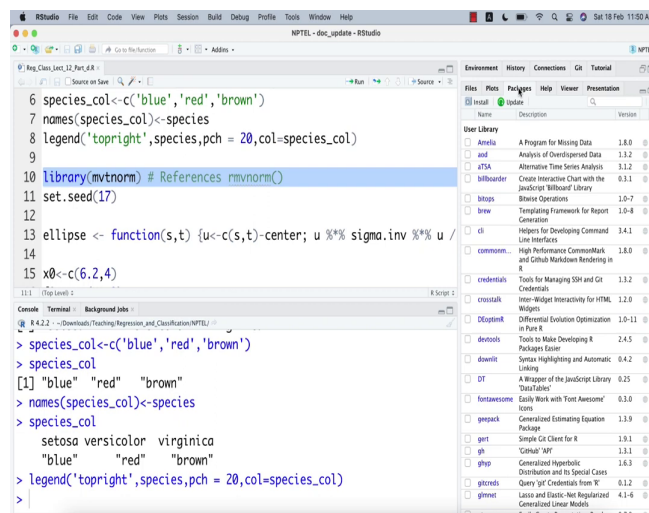
(Refer Slide Time: 00:59)



So, I am going to see the, on the x-axis, I am going to put sepal length. On the y-axis, I am going to put sepal width. So, I am going to name the species. There are 3 kind of species are there.

(Refer Slide Time: 01:18)



And there species color I am going to assign blue, red and brown. And in the species color, I am also giving the name setosa got blue, versicolor got red and virginica got brown.
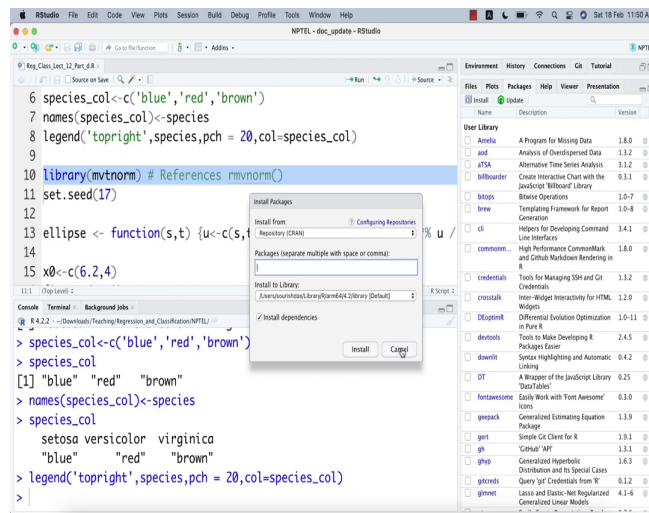
(Refer Slide Time: 01:45)



And then, I am saying that legend put the legend on the top right. So, the setosa blue, versicolor red, virginica got brown. You can try some other color, no problem. Now, I am calling mvtnorm package. If you do not have the package, please install by, here you can go to the package.

(Refer Slide Time: 02:05)



Then install, and write this thing mvtnorm. Make sure your depend install dependencies corrected and then install. That will work.

(Refer Slide Time: 02:20)



I am cancelling it because I already have it here. Setting a seed here. This function, one line function will help us to draw a ellipse or the contours of the Gaussian distribution. So, x naught is the point for its sort of the 6.2 and 4, is the point I am we are going to use as a test point. And then, I am going to create a f naught with the name species. So, f naught is going to be my likelihood or function.

So, i is in species. So, there are 3 species, for each species I am going to work with.

(Refer Slide Time: 03:14)



From the iris data set, from the iris data set, what I am going to do. Whenever, iris species equal to the species name, then you just take the sepal length and sepal width that will be the x, ok. And then, you compute the covariance matrix of these two column and the mean of these two column.

(Refer Slide Time: 03:47)



Now, given this mean and this covariance matrix, you simulate 1000 points and compute the centre of the point and the covariance of the of those points, ok. And then, you calculate the sigma inverse. So, that gives you a; then you use this mu and sigma to calculate the in call that multivariate normal, density put it in the multi dmvnorm is the density of the multivariate normal and plug it in.

This is the likelihood value of the multivariate density for x naught, x naught point, x naught is the point at which we want to estimate. So, what is the likelihood for that point? Then, for n equal to 100, we creating x axis, y axis, and we are saying that, ok, please draw in this line we are 2, 3 lines we are drawing the contour for that particular species, ok.

(Refer Slide Time: 05:04)



So, if let me just run this loop.

(Refer Slide Time: 05:13)



So, here is the viewer, plots. So, here is the plots.

(Refer Slide Time: 05:17)



So, let me just run the points. So, these are the 3 points and 3.

(Refer Slide Time: 05:28)



And then, here is the test point, here is the test point. You see this is the test point. Now, I am going to give equal prior probability. I am not saying, so prior probability for this point belong to setosa, versicolor or virginica is equal. I am not giving any special additional preference to any point. So, the name prior of the, now I just add the names also, and then, I am just computing the Bayes probability, ok.

(Refer Slide Time: 06:03)



And rounding of the Bayes probability up to the 3 decimal process, so here. So, that means, Bayes probability is 0.861. This point belongs to setosa with 80 percent probability to versicolor with 0 point 2.9 percent or about 3 percent probability and 11percent probability that it belongs to virginica. So, with a high probability we will say that most likely this point belongs to setosa.

(Refer Slide Time: 06:40)



Now, we are going to do the, this is for one test point. So, I try to demonstrate how linear discriminant analysis in theory, using one test point. But, most of the time you will see that your data will have many test points and you have to run it like a typical ml setup, machine learning setup.

(Refer Slide Time: 07:16)



So, I am going to take the iris dataset with; here is the iris dataset with all the 150 values, with 50 in each probability. And then, I am going to first Split the data into train and test, ok. Split the data into train and test, ok.

So, I am just, I just played the data. Till the data table if you just say; so, basically in train data there were 33 setosa, 32 virginica and 35 were versicolor. So, they are almost same. So, out of 150 samples, I have taken 100 random samples from the main data and those samples belongs to train datasets. 100 samples belongs to and 50 samples belong to the test dataset. Here in this, I am going to fit LDA, linear discriminant analysis. And I am saying that ok fit all the models.

(Refer Slide Time: 08:50)



So, maybe I will just say head Iris.

(Refer Slide Time: 08:55)



So, maybe first I will fit with length plus sepal width, maybe with petal length, ok with two feature, essentially with two feature. This is my first model. And if we just run this then we can do the same with the predictive class.

And then this is the result 1. And this is the confusion matrix. Essentially, actually this is only the confusion matrix. So, I would rather instead of result I would say it is a confusion matrix. And let me just show you how the confusion matrix looks like, ok. There are two cases where virginica; which were virginica that was classified as versicolor, ok. So, this is the linear discriminant analysis with two feature.

(Refer Slide Time: 10:25)



Now, let me run the entire thing with maybe model 2. And what we can do; head train, I am sorry; Iris.

(Refer Slide Time: 10:41)



```
59                , Iris, prior = c(1,1,1)/3, subset = train)
60 pred_class1<-predict(mod1,test )$class
61 conf1 = cbind.data.frame(pred_class1,Actual_class=test$Sp)
62 table(conf1)
63
64 mod2 <- lda(Sp ~ Sepal.L.+Petal.L.
65                , Iris, prior = c(1,1,1)/3, subset = train)
66 pred_class1<-predict(mod1,test )$class
67 conf1 = cbind.data.frame(pred_class1,Actual_class=test$Sp)
68 table(conf1)
```
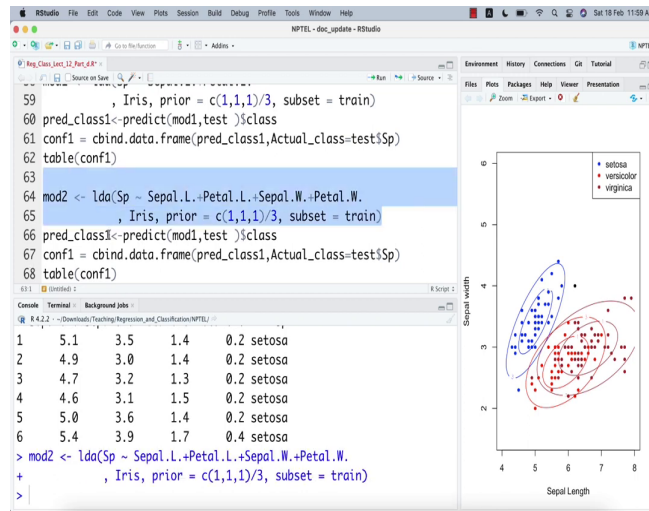
```
> head(Iris)
  Sepal.L. Sepal.W. Petal.L. Petal.W.     Sp
1      5.1      3.5      1.4      0.2 setosa
2      4.9      3.0      1.4      0.2 setosa
3      4.7      3.2      1.3      0.2 setosa
4      4.6      3.1      1.5      0.2 setosa
5      5.0      3.6      1.4      0.2 setosa
6      5.4      3.9      1.7      0.4 setosa
>
```
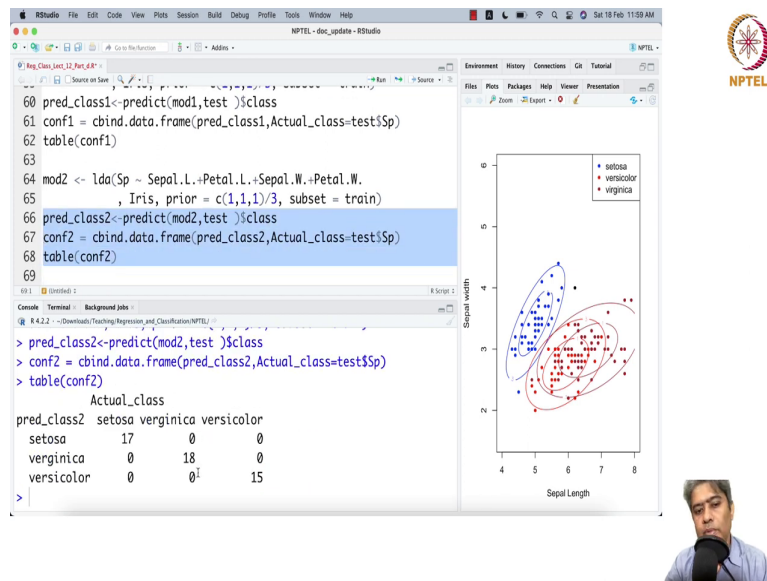
(Refer Slide Time: 10:48)



So, let me now just put the Sepal width and Petal width as well. So, all the features, all 4 features let me put, all the 4 features.

And let me do the prediction. And out of the sample prediction, remember that here is my test data set to do the prediction, ok. And then this is the confusion matrix. So, let us run this confusion matrix, now it is being corrected. So, when we put the more feature, naturally the classification becomes better. So, and it is in the out of the sample classification it got better.

So, this shows that if you have more feature, more typically it helps you. But at the same time, you have to be careful about the over fitting that if you put too many features too many engineered features, it might, you might end up in a wrong place with lot of over fitting. That means in your test accuracy and out of the sample accuracy. And in sample accuracy will be very different.

At the same time, remember that you have to be this particular data set iris data set is actually a very toy data set. Going forward we will try with more real life new days data set, new age

data set where it will be much much difficult to classify the target variable. With this, thank you very much. Wish you a happy weekend. See you in the next week with new lecture, new topic.

Thank you very much.