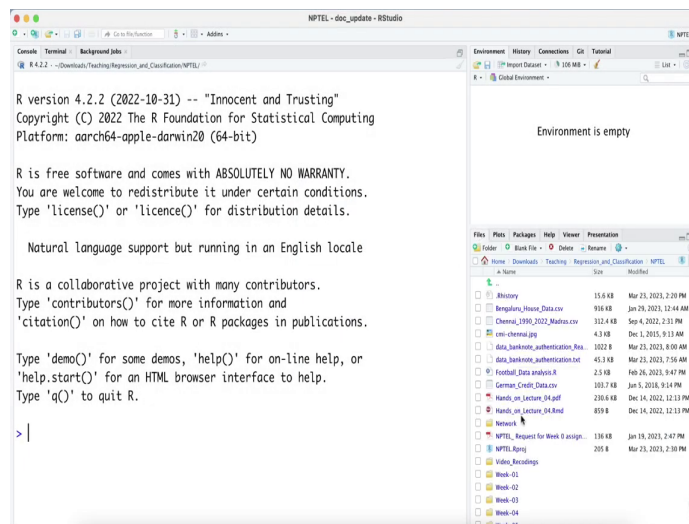**Predictive Analytics - Regression and Classification**
**Prof. Sourish Das**
**Department of Mathematics**
**Chennai Mathematical Institute**

**Lecture - 58**
**Hands on with R: Classify fake bank note with GLM**

Hello all. Welcome to the part B of lecture 19. In this part we will do Hands on with R in this hands on video. We are going to do work on another data analysis where let me just first open the R.
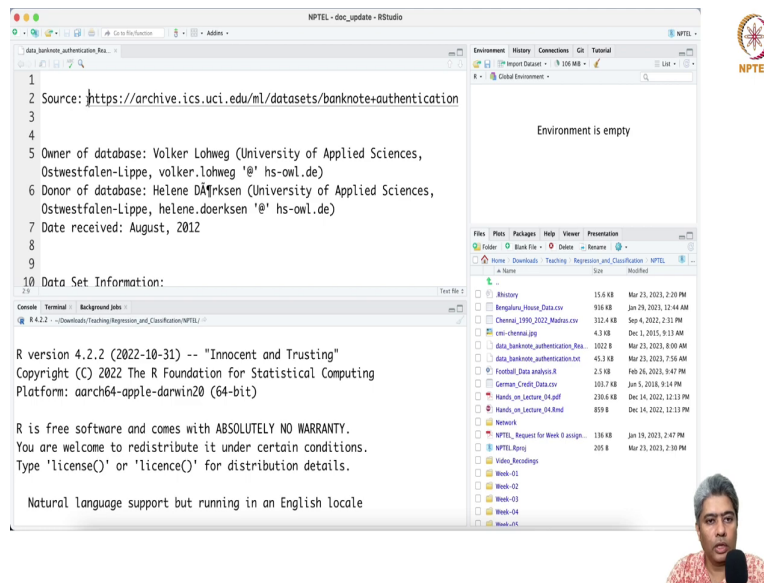
(Refer Slide Time: 00:44)

(Refer Slide Time: 00:58)



So, what I am going to do here we have the from the UCI machine learning repository. If we just go there and let me just open my net, yeah ok.

(Refer Slide Time: 01:18)

(Refer Slide Time: 01:29)

(Refer Slide Time: 01:34)



So, this is the UCI machine learning data repository. This is called bank note authentication data set. So, this is from Denmark. It says basically data were extracted from images that were taken from genuine and forged bank not like specimens. For digitization and industrial camera usually used for print inspection was used.

The final images have a 400 by 400 pixel due to object lens distance to the investigated object grayscale pictures with resolution of about 600 dpi where gained. Wavelet Transformation tool were used to extract the features from the images.
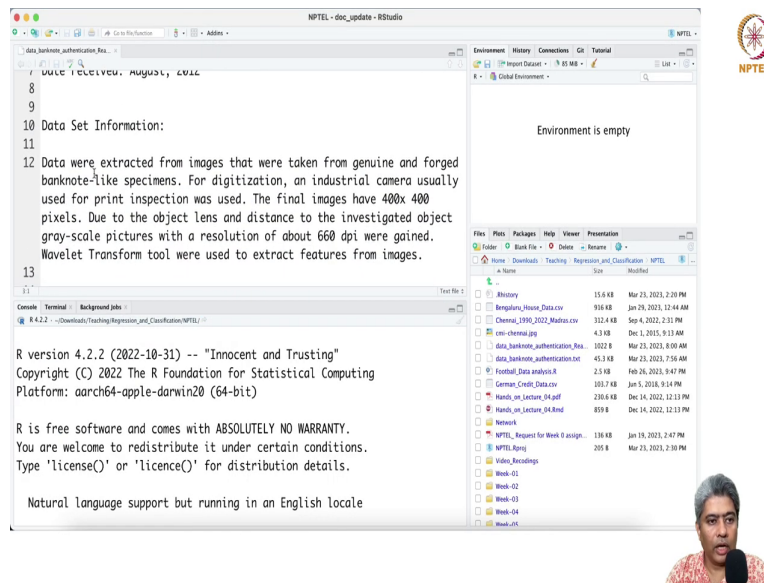
(Refer Slide Time: 02:15)



And these variants of wavelet transformations, skewness of the wavelet transformations, curtosis of the wavelet transformations, entropy of images and the classification manually done classification 0 or 1 that was done.

(Refer Slide Time: 02:40)



So, now let us start analysis. So, data equal to we will do read dot csv read dot csv. Let me just go quickly to the bank note and all downloads and classification. So, from so, this is the data set name. I have shared this data with in the swarm portal.

(Refer Slide Time: 03:25)

(Refer Slide Time: 03:28)



So, if we just read this ok. So, it does not have any header because looks like this cannot be the name. So, it does not have any header.

So, we have to say header equal to header equal to false header equal to false and now we have the data correctly read, alright.

Now, what we have to do? We have to put the colnames data equal to c. First, let us take these name variants of wavelet. So, first one will be var of wavelet, second will be the skewness, curtosis and entropy, ok.

Skewness of wavelet maybe I will just say skew. I do not want to put very large name. Skewness, curtosis of wavelets and "entropy"; and finally, last one is the class which is forged or not, right. So, let me just run this. So, now we have the these names are being given ok. So, maybe I will just say wave. Yeah. And now I think this is fine, right, ok.

Next, I will just have a look that how many forged cases and how many original note were there. So, out of 13,722 13,000, no, sorry, 1372 observations, 610 of them are forged and 762 of them are not forged. They are true original bank note. So, what I am going to do?

I am going to first plot data, dollar, variance of wave, comma, data, dollar skewness of wavelength maybe pch equal to 20. Now, what I am going to do is something you will I think

you will have fun with it, data, dollar, you take the forge or not, see there is 0 or 1 is there. So, what I am going to do? I am going to add 1 to it, ok. I am going to add 1 to it.

(Refer Slide Time: 07:10)

(Refer Slide Time: 07:15)



Now, what happens if I add 1 to it? What is happening? You actually is very simple. I am just adding either 1 or 0. It was 0 or 1. So, all 0s have become 1 and 1s have become 2. Now, as a color, 1 color equal to 1 means it is it will be black and color equal to 2 means it will be red, ok.

So, this is a this is I am taking making forged or not as the color that I am going to use xlab is let me just take simply "variance of Wavelet", ok. And ylab, ylab equal to this is simply "skewness of Wavelet" transformed image, ok.

Let me just make this plot. Yeah, let me just zoom it. So, this is the data. So, now we see this data is very peculiar. The data has a very peculiar shape. You can see there is a very peculiar shape and there is bit of a overlap between the two kind of things. So, one is was forged, which was marked as red or these are like these are the specimens or points which have which represents forged note and the black ones are the one, which represent the original note.

So, there is we see there is a bit of a overlapping is happening, but mostly they are kind of well separated with slight overlapping between the 2. So, can we use wavelet variance and skewness of the wavelet these two feature to do the classification.

(Refer Slide Time: 09:19)



So, what we will do we will first we will set a seed say sample first we need a set a seed. So, we need a number 1 is to say 1000. So, we will take 81. So, 81 as sitting our seed set dot seed ok, alright. And then n is the nrow of data and m is the ceiling of n times 0.7.

So, I am taking 70 percent for training data set and rest of them has test idx equal to sort, sample, 1, is to n comma m comma replace equal to false, ok. And then so, let me just have this. So, this idx is 961 cases have come and then now df train equal to data idx comma and df test equal to test equal to minus idx. So, let me just run the whole thing, yeah.

So, first thing I am going to do I am going to fit a simple model with wavelet variance of wavelet and the skewness of wavelet as the you know feature of simple logistic regression fit. First model will going to be glm, ok glm. So, what we have is forged or not is as this comma variance of wavelet plus skewness of wavelet.

These two data equal to data equal to df train and then family equal to binomial link equal to "logit", alright. Now, this is what I am going to do. I am going to run this model. Now, if you run this summary fit1.

(Refer Slide Time: 12:34)



So, what you have it is saying that both variance of wavelet and the skewness of the wavelet is going to have a very strong accuracy sorry strong effect on the both variance of wavelet and the skewness of the wavelet is going to have a on the whether note will be forged or not, ok.

So, let us take that as and then now what we are going to do we are going to make the prediction on the test data set, ok. df test df test dollar probability1 equal to predict fit1, newdata equal to df test and then type equal to "response", ok.

(Refer Slide Time: 13:54)



So, now if you see you see this, alright. So, now you can see that these values have come 0 and all this probability of 0.24.

(Refer Slide Time: 14:17)



So, let us see if it is. So, mostly they are close to 0 some 0.1, 0.2 and if we just go down to 1 and you will see that they are reasonably close to 0.9, 0.8, 0.9, 0.8, right. Here is a one point where you got a you know 0.26, but the probability is less, but you know actually it is a forged note. So, there will be looks like there could be some misclassification possible. So, what we are going to do? We are going to make the prediction. These are the probability of this particular node to be forged note.

(Refer Slide Time: 14:51)



Now, I am going to make the prediction df testdollarprob1 equal to df testdollarprob1. If it is 0 greater than 0.5, we will call it 1. And if it is less than equal to 0.5, we will call it 0. So; that means, not forged ok, alright.

Now, what we are going to do? So, now, if I just look into the df test. So, now, they are all. So, here are there are some mis-classification we can see, but mostly they are agreeing.
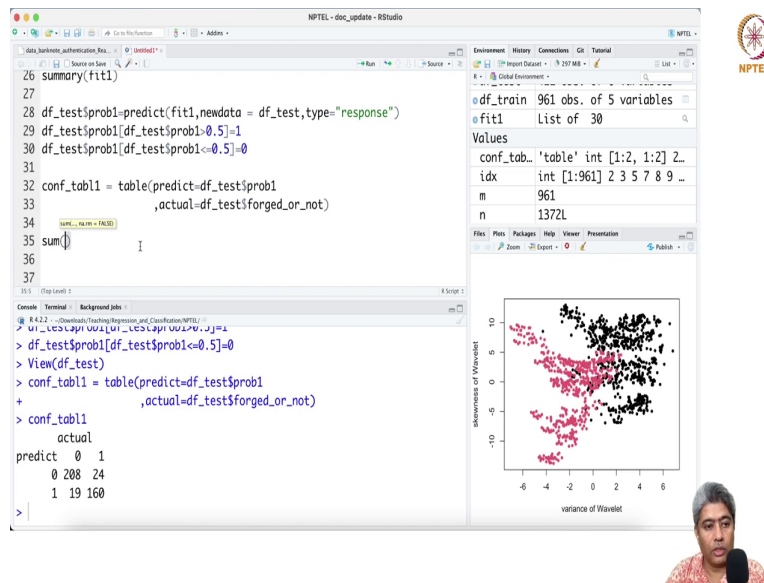
(Refer Slide Time: 15:45)



So, now we are going to I am going to calculate a confusion table, ok. So, table predict equal to df testdollarprob1 comma actual equal to df test dollar forged or not, ok. Now, if I just run the confusion table. So, this is my confusion table actually predicted 0 and actually 0 was 208 cases and in is the test data set. Similarly, predicted 1 and actually it was 1 was 160 cases.

There are 24 cases where it was predicted it is not forged, but it is actually forged and there are 19 cases where it was predicted forged, but it was not actually forged. So, there are possibility of and let us just check the accuracy.

(Refer Slide Time: 17:06)



So, first I am going to take the diagonal of the confidence table. These are the correct predictions all correct prediction take the sum of that; sum of that. These are the all 368 are the total correct prediction and total number of all test is. Actually, just I can take nrow of df test is 411 of them there are 411 of them if I just make 100 out of this so, 89.53. So, this is my accuracy 89.53.

Now, if I make, I will make a second model fit2. What I will do is essentially I will just update fit 1 with same model plus few more engineered feature. So, I will just take variance of wave and square them plus a skewness of wave lengths and square them and finally, plus variance of wave times skewness of wave, ok. If I just run this and summary let me run the summary of fit2.

(Refer Slide Time: 19:20)



So, looks like variance of wavelets and skewness of wavelet ok was. So, initially variance of wavelet ok if I just run the summary of fit1 we see that both wavelet both variance and skewness of the wavelet is as effect, but if I run on the forged or not whereas, if I run summary of fit2 where, it shows basically variance of wavelet does have a effect skewness does not, but the quadratic effect of skewness does have a effect and their interaction does have a effect.

So, this is very interesting phenomena, but let us see what is the whether adding this feature help me increasing the out of the sample accuracy or not. So, that will be my that will be the real test case.

(Refer Slide Time: 20:33)



So, what I will do? I will just take this guy copy and paste it here, but here instead of probability1 I have to take probability2 and then instead of fit2 I will just do fit1 I will just do fit2 it will add a new column it should add a new column the probability2, ok.

(Refer Slide Time: 20:47)



And then on the probability 2 I will make the prediction if it is more than 5, 0.5 I will say 1 otherwise I will say it is 0. So, now they are making 0 1, ok. Here is a misclassification; here was misclassification it was corrected here is a misclassification here is misclassification.

(Refer Slide Time: 21:26)



So, we will see and we can just calculate a confusion table simply second confusion table with probability2 and what is the situation of this is the second confusion table for this with the second model and if we just compute the accuracy simple accuracy this is 91.97. So, 92 percent accuracy we are getting whereas, for from this first model we are getting 89.5 percent.

(Refer Slide Time: 22:02)

(Refer Slide Time: 22:23)



So, nearly 2.5 percent gain that we are getting, ok. So, this could be accuracy1 and this could be accuracy2 and accuracy2 minus accuracy1. So, 2.4 percent gain that we are getting because of the because of engineering proper feature engineering. Now, let us see how the geometry of the predictive space has changed using because of adding the feature engineer engineered feature.

So, I will do some visualisation to understand the effect of feature engineering, ok. Let us try to understand that. So, first I am going to do test data I am going to define a test data or dummy variable and dummy variable data dot frame and with that I will going to see matrix NA comma nrow equal to 1 ncol equal to 1 or sorry ncol equal to 2.

(Refer Slide Time: 23:50)



So, test data and this is the I am just creating a simply place holder, ok. Just a place holder and then I am going to add a colnames of a test data equal to colnames equal to dummy variable equals to this 2, ok. These are the then variance of wave, alright. So, variance now what I am going to do? So, variance of wavelet as a range of between see it is between minus 7 and 7 and skewness of wavelet is somewhere between maybe minus 12 to 12, ok. So, I am going to create this variance of wave define minus 7 to 7 by 0.1.

(Refer Slide Time: 25:12)



So, basically if I just do that now you can see that it just creating values between minus 7 to 7 with a difference of 0.1 ok, it just created these values.

(Refer Slide Time: 25:28)

(Refer Slide Time: 25:44)



Similarly, if I just say skewness of wave is going to take sequence of value between minus 11 to 11 by 0.1 difference this is also going to create values between minus 11 and 11 with a difference of 0.1, ok equal distance of 0.1. So, I am getting grid I just created grid.

(Refer Slide Time: 26:01)



Now, I am going to create a joint grid; joint grid how that is how. So, for i in 1 i in variance of wave for j in skewness of wave and then dum variable equal to c of whatever i and j values are there in the thing and then what I am going to do? Test data equal to rbind data dot frame test data comma dummy variable, alright.

(Refer Slide Time: 27:13)



So, let me just run it, ok. So, it has created test data with 31,162 observations if I just create that you will see that it has created all a complete grid of two values.

(Refer Slide Time: 27:27)



And now what I am going to do I am going to create a let me just put a na dot omit dot omit. So, there will be the first row will be deleted and now you can see all the values are being created.

(Refer Slide Time: 27:48)



Now, essentially what I am going to do is I am going to create for these all possible values I am going to create probability1 equal to predict fit from the 1st model newdata equal to test data and type equal to "response", ok.

Now, what I am going to do is test data from the test data I am going to do a prediction and here is the whole thing the prediction is basically probability1 if it is greater than probability1 greater than 0.5 then the prediction will be "pink" basically it is a forged; that means, its forged, ok and if it is less than equal to 5 we will use grey; that means, it is not "grey" color I am just going to make this things, ok. Now, if you just go there you see pink and grey.

(Refer Slide Time: 29:37)

(Refer Slide Time: 29:44)



So, now what I am going to do? I am going to make a plot I am going to make a plot that test data dollar variance of wave test datadollarskewness of wave and then pch equal to 20 color equal to test data dollar my prediction color whatever the prediction color. Now from the above plot so, I can I will just take the lab labels xlabels and ylabels, ok.

So, that is the my predicted area. So, later the one which so, any point that will fall in this area will be called forged note and any points or any notes whose variance of wavelet and the skewness of the wavelet falls in this area will be called not forged or original note, ok. And then what I will do? I will I can just put the points essentially df testdollarvariance of wave comma df testdollarskewness of wave pch equal to 20 and color equal to df test dollar the forged or not plus 1.

So, I have to just take this and then alright. So, if I just run this. So, now, if you see so, these are the points. So, let me just. So, clearly you can yeah I think this is now probably better right. Now, this looks good.

So, what we are seeing here that these are the points you can see there are some misclassification point the points which were definitely where will be misclassified and then there are like you know forged note, but they will be called as a original note and then there will be notes which are original, but they will be called as forged.

So, there is a effectively bit of a difficult this is sort of a difficult zone, ok. This is definitely a bit of a difficult zone, alright. Now, what I am going to do I am going to produce the same

plot, but for the second model remember that second model were second model were developed using the engineered features with quadratic terms, alright.

(Refer Slide Time: 33:44)



So, I am going to copy this thing. So, probability2 with second model probability2 probability2 prediction2 probability2 prediction2 probability2 let me just run this.
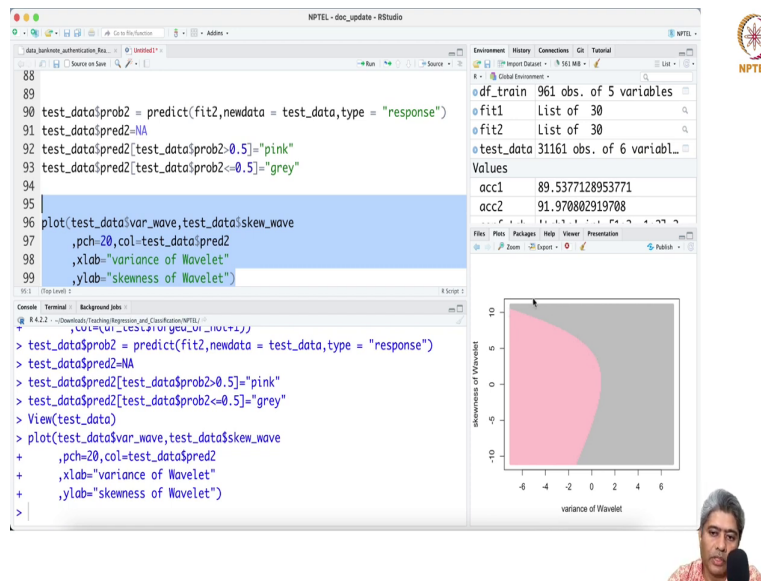
So, now if I have the test data, I have now got the second probability and their corresponding predictions as well.

(Refer Slide Time: 34:27)



So, now if I just run this, but the color will be instead of prediction1 the prediction2. Now if I just do that. So, now, you can see this is a completely different geometry a quadratic geometry sort of coming up nicely, ok sort of a ellipse or kind of thing behavior is coming up, which is a very interesting and now if we just put the points here, right.

(Refer Slide Time: 35:00)

(Refer Slide Time: 35:06)



So, put the points here then we can see that you know there is a it is trying to capture these points whereas. So, let me just you know little bit increase this guy. So, it will be let, yeah.

(Refer Slide Time: 35:30)



So, this is the quadratic behavior that we are seeing here and whereas, if you just this is the simple without any feature engineering, we will get a simple linear you know discriminator whereas, here we are getting a sort of a nice non-parametric quadratic discriminator. Here you can see these points are getting misclassified right, these points are getting misclassified whereas, here they are not getting misclassified this point is getting misclassified.
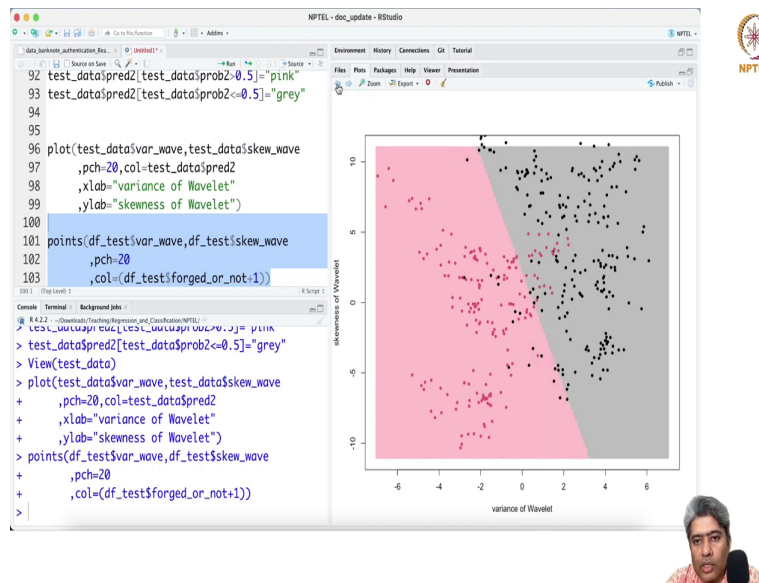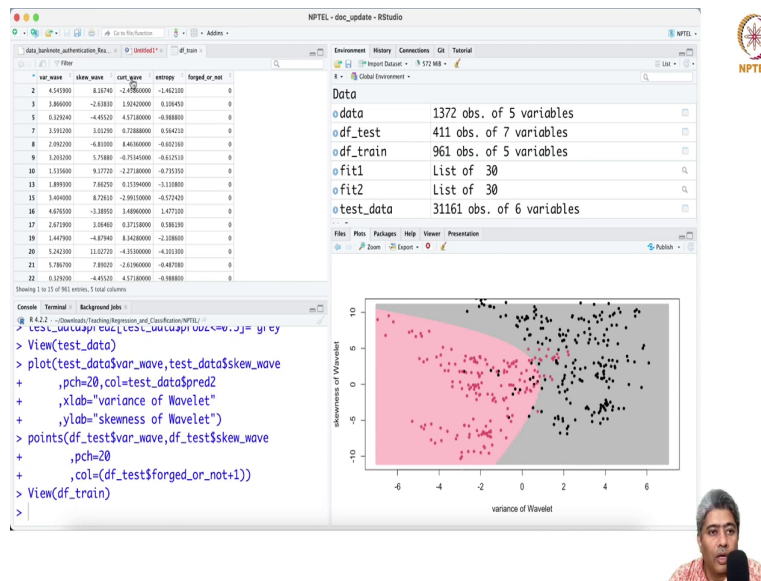
But it is better than too many and similarly I think if this point is getting misclassified here whereas, sorry about that, this is better this point are not getting misclassified. So, few points which are nicely not getting misclassified and trying to capture this and overall amount of the sample accuracy is going up.

(Refer Slide Time: 36:51)



So, feature engineering typically helps remember that I have not tried the cubic and the other higher order transformations, in addition we have few more data for example, we have curtosis and entropy, which I do not have I have not used. So, my recommendation is you guys use that and check whether it helps you to improve the accuracy.

So, I will stop here and in the next video, I will try to do predictive analysis with a new data set with a new problem. I hope you are enjoying these hands-on real data analysis. Take care, bye.