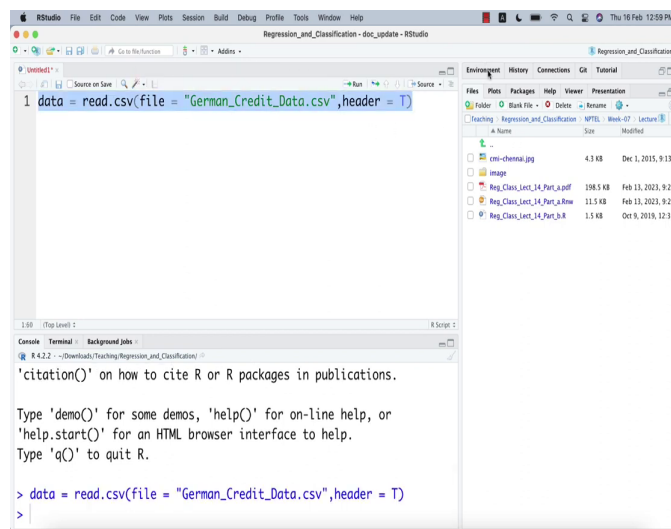**Predictive Analytics - Regression and Classification**
**Prof. Sourish Das**
**Department of Mathematics**
**Chennai Mathematical Institute**

**Lecture - 38**
**Hands on with R for Logistic Regression**

Hello all, welcome back to the part C of lecture 11. In this part, we are going to do some Hands-on.

(Refer Slide Time: 00:23)



So, we are going to call the and this hands-on will be based on R. So, first we are going to read the CSV file. So, we are going to read the German credit code data. So, I have already shared this on the NPTEL platform, Credit also if you do search German credit code data on

internet, you will get this data same data set its a slightly old, but reasonably good data to get your hands dirty with, you know, the first logistic regression model kind of thing.

(Refer Slide Time: 01:36)

(Refer Slide Time: 01:54)



So, yeah, if you just come see there are about 1000 values and there are some status are given and if you look into the status, if you go to and then the target here is good or bad, good loan or bad loan and then essentially I, we have to, it is recorded as 1 or 2.

(Refer Slide Time: 02:02)

(Refer Slide Time: 02:29)



So, what we will do is data dollar Good Bad equal to, if we just say data dollar Good Bad minus 1. So, that will give us, you see good or bad, if it is 1, then it is bad, 0 means good. So, we want to predict accordingly.

(Refer Slide Time: 02:40)

(Refer Slide Time: 03:07)



Now, if we just do summary, summary of good or bad. So, about 30 percent is bad actually. So, just a minute, I think if I just do table instead of summary, the if I just do a table. Yeah. So, 300 out of 1000 cases, they have 300, they have given as bad cases and 700 given as a good cases. Now, typically in real life, you will not have default more than 2, 3 percent of the entire credit customers.
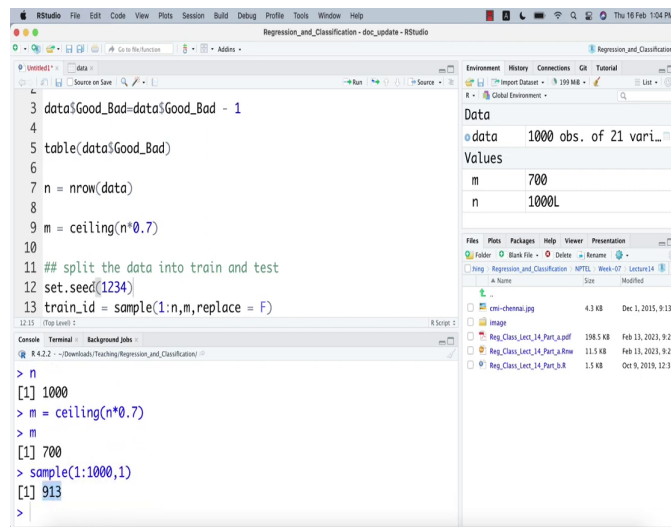
So, this is slightly a artificial in that sense because you cannot, if you are, if your 30 percent of customer is this defaulting, Then you will not be able to you know survive. But anyway, but if you have in your data set, if you have only 2 percent ones and 98 percent zeros, those kind of data sets are called imbalance data sets and in those kind of data sets doing prediction is extremely difficult.

And because most of the time, what happens that any standard model will just say everybody is kind of, you know, good, then no matter what, what is your, what happens 98 percent of the time you are correct. So, and in only 2 percent time you are wrong. So, a naive model will do correct, will be correct 98 percent of time. So, in real.

So, those kind of naive model, we have to be very careful about. So, what I will do first, I will just say number of rows in the data is 1000, I believe yeah. And then I will take m is the testing number of testing, like, you know, n into 0.7, maybe 70 percent of the data we will consider for the testing.

So, 700 out of 1000 700 data points will be considering for testing. So, for we will do a split the data into train and test. First, what we will do? we will choose the train id's sample 1, is to, n out m many samples you will we will draw and replace equal to false. Ok, we have to set up a seed.

Set dot, seed is 1234. So, the advantage of using seed is if you set up a seed and if you use exactly this number 1234 actually you can do any number, I can just draw a sample between 1, is to 1000, one sample and 913, I am let me just use 913 as my seed.

(Refer Slide Time: 06:34)



Now, that means, if you use this seed number as your seed, then what will happen is you, you my result and your result will exactly match ok. So, this ensures because we are going to split the data set randomly 70 percent data we will take as a training data, ok. So, data underscore train equal to what I am going to do data. Just I will provide the training id comma maybe I will just do a sorting here ok.

Let me just take this and then I will just data train and data test equal to what I will do. I will just copy this and instead of train id, I will just say minus train id. So, just take all the ids, but do not take the training ids.
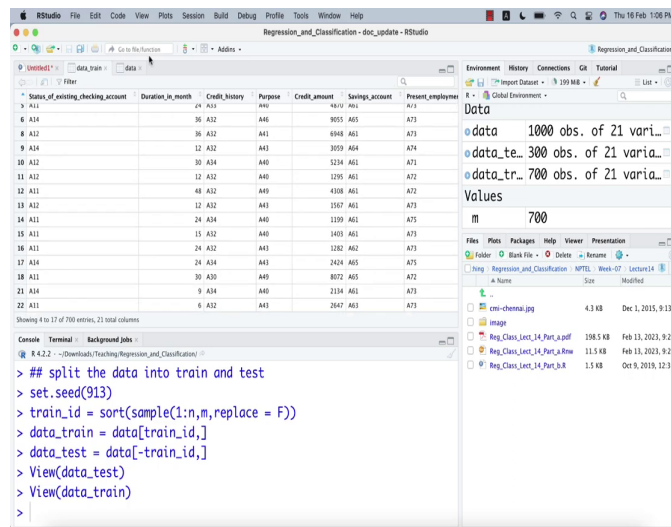
(Refer Slide Time: 07:40)

(Refer Slide Time: 07:57)



So, that will automatically give me the testing data. And so, row number 1, 7, 19, 20, 26, this came to the test data. And if you go to the training data, 2, 3, 4, 5, 6, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, then 19 missing 21. So, those are the train data, ok. So, now we have splitted the data. So, first we will fit a logistic regression using built-in function in R ok.
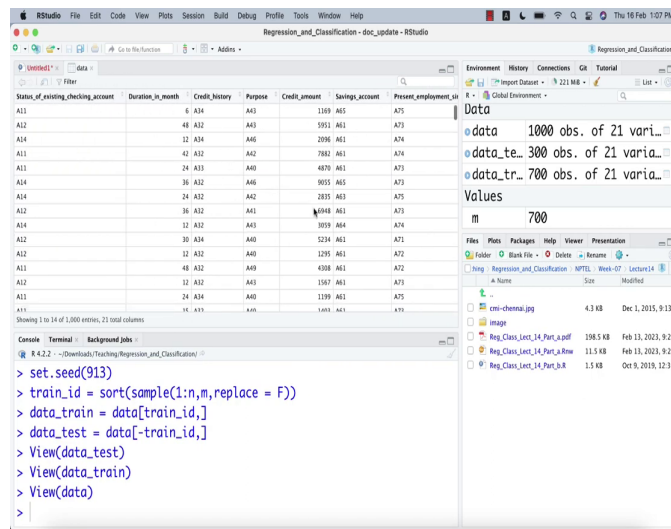
(Refer Slide Time: 08:19)



First we will fit logistic regression using built-in function built-in function in R ok. Now, first fit glm. So, from the stats package, you extract the glm function, you call the glm function. So, you just say Good Bad and dollar.

(Refer Slide Time: 09:14)



And I am going to take two in this data. There are most of these data are categorical, but I just want to keep it little simple. So, that, you know, you understand it more care more easily.

(Refer Slide Time: 09:40)

(Refer Slide Time: 09:44)



The duration in month is a sort of a numeric value, you know, yeah, credit amount and duration in month. These are the 2 columns, which are in numeric. There are some other numeric also, installment rate, present, residence, since age. So, these also can be, I can take, but just to be, my purpose of this demo is to give you how the logistic regression works, ok.

You can add as many predictors as you want. There are about 20, 19 or 20 predictors are there in this data set. But you can add as many predictors as you want. My purpose is to just give you how the whole logistic regression works in R.

And what I am going to do, I am just going to take 2 predictor, which are numeric predictor. And what I am going to do, I am going to just scale. This scale always helps Duration in month and plus scale Credit amount credit amount and then data equal to data train.

And now here is the important part that you have to mention family, family equal to binomial, because remember that it is a binary classification model. So, it raised into binomial I have to save it somewhere. Let me just call it. So, I will just call it maybe, MLE logistic regression logistic regression, ok binomial. So, since it is a binary classification problem, in the binary classification problem, the model automatically raise up what it comes is a binomial model or Bernoulli model.

(Refer Slide Time: 12:07)

(Refer Slide Time: 12:35)



So, and we have to give a link, link equal to, we will say, logit. You can use probit also. You can use probit also. And let me run. Let me run. It should be stats, not stat. Yeah, ok no error means probably it has run well.

(Refer Slide Time: 12:49)

(Refer Slide Time: 12:52)



Now, I am going to say summary, fit glm ok. So, it what it says that duration in month and credit amount, both of duration in month has a statistically significant m effect in the, in the whether a probe, customer is going to be a good customer or bad customer.

And turns out credit amount is not that important. You see, p value is very large, but credit amount is important. So, this is the estimate, the coefficient. These are the standard error and z value. So, you can have these, you can create these information's. And now what I am going to do is now we are going to pretend that as if glm function in our does not exist and we want to code it from scratch ok. So, we are going to write down the negative log likelihood function.

And then we will use the optimization subroutine technique and we will find what estimates of the coefficients are. And what these estimates are going to be match going to match with built in glm function ok. So, we will first write negative log likelihood function ok, alright.
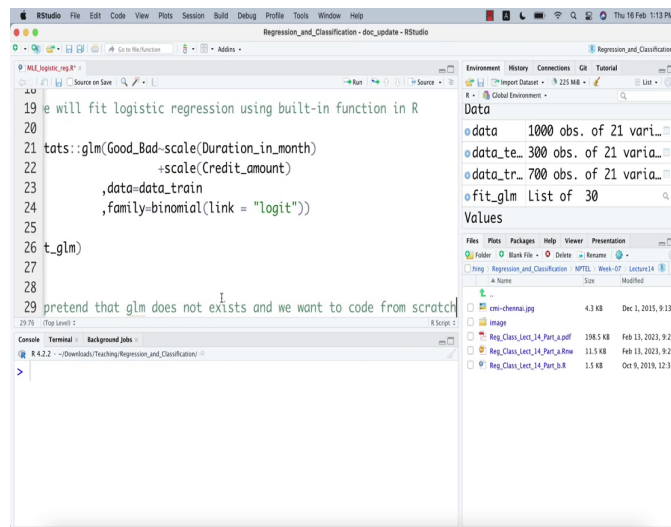
(Refer Slide Time: 14:21)

(Refer Slide Time: 14:54)



Let me first, let me just write a point here that: Now we will we pretend that glm does not exist exists and we want to code from scratch ok.

Now, we want to code from scratch. Now, first, how do you do that? Negative log likelihood you have to write the negative look likelihood function. So, there will be all the statement will be inside and then everything is first dataset name we have to give then beta values we have to give then y name we have to give and the x names we have to give correct, alright. So, the first from the we have to extract y from the dsn y name and keep it in x and then x equal to scale dsn x name.
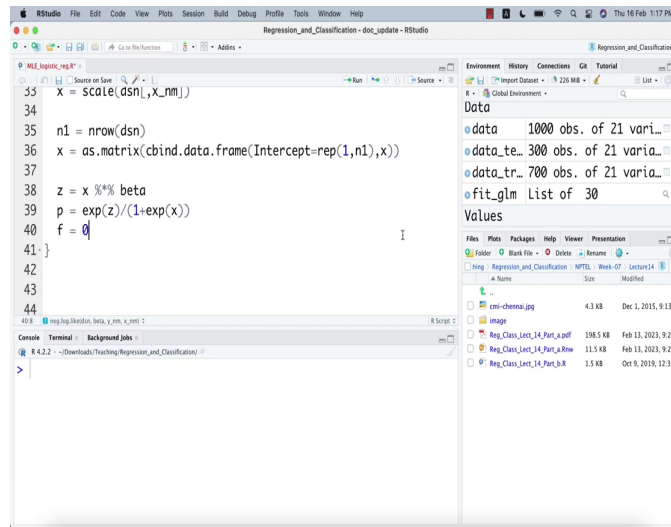
So, the way I am doing it at this moment it will not work for categorical variable. So, it will work only for the numerical variable when writing the full blown code is not my purpose. And my purpose is to give you a demo that how if require you can write your own model you can develop your own model write down the likelihood function, pass it through the

optimization sub routine and feed your own model and hey, here is your own model isn't it fun?

So, what I am going to do is first is n1 is nrow of dsn and then I am going to define it as remember that most likely dsn is going to be a data frame not most likely we are going to pass it as a data frame, but we need to do a matrix operation. Since we need to do a matrix operation what we need to do we have to call it as matrix this x as dot matrix. And we have to do this cbind dot data dot frame with x and before that Intercept equal to replicate of 1 comma n1.

So, we have to create a x matrix with all the x columns that has a name of the x columns with all the predictors name and a column, which will be intercept right. Now then once I have the x what I have to do I have to calculate for z remember that in the theory part we have a z value latent variable z. So, this is the latent variable I am going to introduce here. Percentage star percentage beta and then p equal to e to the power z divided by 1 plus e to the power z.
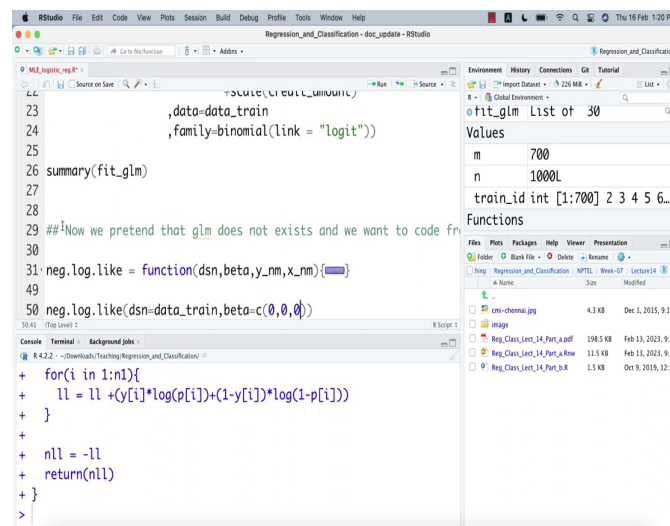
(Refer Slide Time: 18:50)

Now, I have the p wonderful now what I have to do is f equal to 0 I have to calculate the log likelihood actually I should say log likelihood ll equal to 0. For I could have done it slightly better way in a some kind of you know more computed efficient way, but my purpose is to show you not a write a write an extremely efficient way of writing the code when my purpose to show you how the concept of writing negative log likelihood.

And because the best thing is if I could have write this three lines of code in a vectorized form, but instead of writing it in a vectorized form I am purposefully writing it as in a loop so, that you can understand what I am doing. So, plus now y i times log p i plus 1 minus y i into log of 1 minus p i ok.

Now, here is a here is a task for you ok. Write this code three lines of code without for loop, can you do that and check if do the, if you do that whether your overall optimization time is

going to be reduced or not. So, let me just let me just come finally, write negative log likelihood will be just minus of log likelihood. So, lll this is my log likelihood and just and then return negative log likelihood ok.
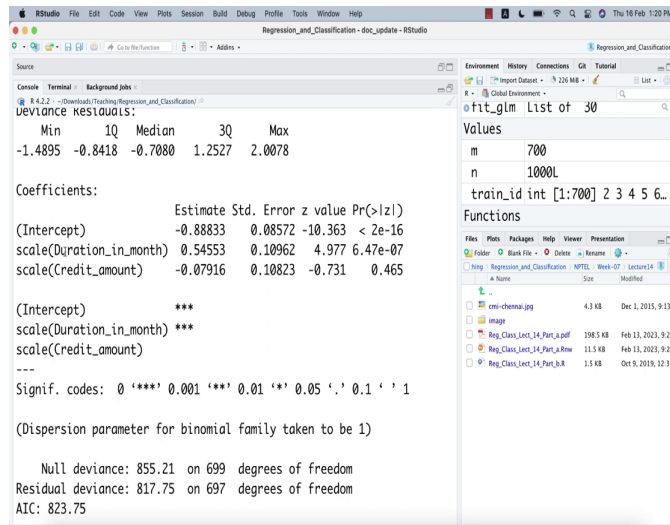
(Refer Slide Time: 21:34)



Now, once I have written it once I have written the negative log likelihood at least we have to check whether the value the function works for some data set or not. So, some value some initial value dsn will be data underscore train then beta will be 0 comma 0 comma 0. So, why is 0 comma 0 comma 0? Because you see I have if I have in the model that I am trying to fit this model I am going to fit.

(Refer Slide Time: 22:08)



```
Deviance Residuals:
    Min       1Q    Median       3Q       Max
-1.4895  -0.8418  -0.7080   1.2527    2.0078

Coefficients:
                          Estimate Std. Error z value Pr(>|z|)
(Intercept)               -0.88833    0.08572 -10.363  < 2e-16
scale(Duration_in_month)   0.54553    0.10962   4.977 6.47e-07
scale(Credit_amount)      -0.07916    0.10823  -0.731    0.465

(Intercept)              ***
scale(Duration_in_month) ***
scale(Credit_amount)
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 855.21  on 699  degrees of freedom
Residual deviance: 817.75  on 697  degrees of freedom
AIC: 823.75
```
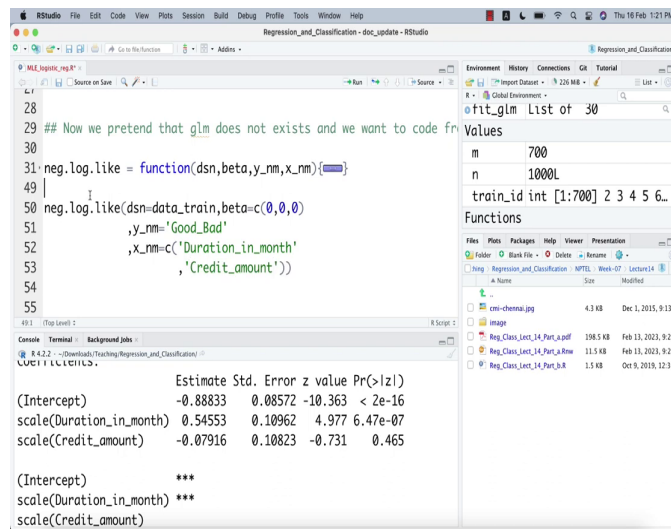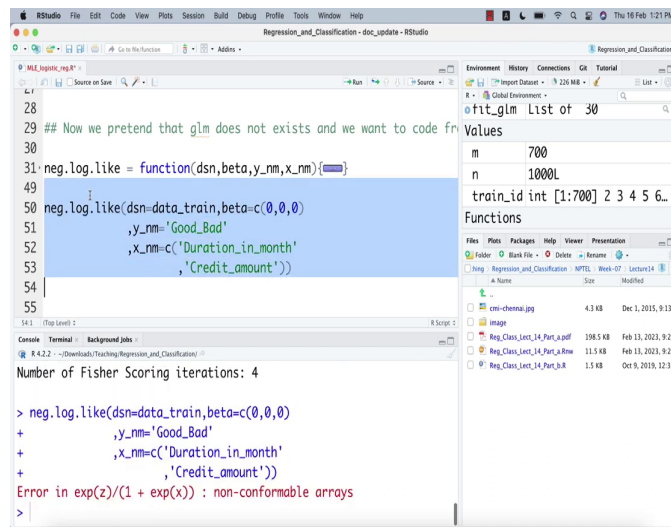
So, intercept one coefficient for Duration in month another coefficient Credit in amount. So, these are three coefficients for all three coefficients I am going to give initial value of 0. And then I have to give y name y name y name is in a vector form sorry y name is good or bad Good Bad and x name actually I can just copy these names Duration in month. Because that is the number exactly what we have and Credit amount ok.

(Refer Slide Time: 23:09)



So, if I just run this these are the initial values this are called non-conformable let me see of course, because it should be z. So, there was a mistake alright.

(Refer Slide Time: 23:27)

Let me run it once more yup. So, in the first initial value the negative log likelihood function runs and it gives me a value of the log negative log likelihood now I am going to initiate the initial value of beta I have to give initial value as 0, 0, 0 I have given here maybe I will just run it here and I will do it there.

(Refer Slide Time: 24:00)

(Refer Slide Time: 24:27)



And if I just run it yeah it is working perfectly fine and then I am going to write optimization sub routine to get the maximum likelihood estimate write optimization or just call optimization sub routine optimization to estimate MLE of beta ok. So, what I will do, fit lr logistic regression with optim and from the stats package you just call optim ok the parameter first is you have to give the parameter which is initialize beta dot init.

Then you have to give the objective function on what is your objective function? This is your objective function negative log likelihood function. So, you have to give the objective function then you have to give some additional in the objective function you have dsn data dot train. Then in the objective function you have this y name and x name that you have to provide ok.

(Refer Slide Time: 25:55)



So, and if you run this, I have made again a mistake it should be stats not stata and it has run and without any error. So, this is a good news.

(Refer Slide Time: 26:26)



And let me run from this let me call the parameter estimated parameter value. So, these are my estimated parameter values. So, these are my parameter values from manually estimated whereas, when we use the summary dot fit or maybe coefficients we just call coefficient of fit glm this is looks like these are exactly not exactly.

But at least up to three decimal places they are matching looks like 0791, 0791. So, up to 3 to 4 decimal places it is matching three 3 to 4 decimal places it is matching. There are ways to match it if I increase the tolerance level and all these things then it will match maybe up to a great extent.

So, what I will do if you. In fact, I can just say cbind and to just compare we can see. So, first one is from optim and second one is from the built in yeah. So, this is the when we use optim this is the value and when we use the built in r function glm this is what happens. So, they are almost matching up to the three decimal places and then they are deviating this one has matched up to four decimal places three to four decimal places they are matching.

If we increase the tolerance level probably, we will match even more. So, now, we have the let us write the beta hat so, since ok. So, beta hat equal to fit lr par and then x test x test we have to take from we have to scale it remember that we have to scale it data test comma these two value.

(Refer Slide Time: 28:41)

So, we have to take the test ok we got the test values and then what I am going to do is we are going to calculate this first the latent variable z hat which is x test percentage star percentage beta hat. Sorry I have to call this guy and then yeah it is not confirmable. So, of course, because I have to do few more work in row of x test and then I have to call x test equal to you see I have done this line right.

So, I have to do the same kind of lining arrangement here also because that is not x test directly is not my you know my call. So, let me just run this and now probably it will be fine, yes now it is fine now it is fine. And now once I have the set test let me just have p hat p hat equals to e to the power z hat divided by 1 plus e to the power z hat ok. So, this is my p hat. Now, now in the data test data test I am going to p hat equals to p hat this is going to be a new column in the data test. See I have created a new column here alright.
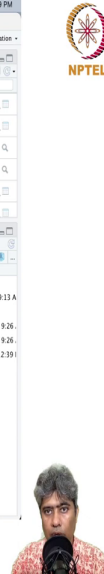
(Refer Slide Time: 31:17)



Next what I am going to do I am not sure why it shows a p hat column 1 ok, we will see what can be done alright. Next is we have to predict.

(Refer Slide Time: 32:08)

(Refer Slide Time: 32:24)

(Refer Slide Time: 32:32)

(Refer Slide Time: 32:38)

So, data test dollar Good Bad, but we have to do first prediction the prediction equal to first I am saying 0 and then first I am taking all the values of sort of a 0. Actually, what I can do? I can just take this yeah, I will see what (Refer Time: 32:44) can do and then I can just say that if p hat is greater than 1 greater than half. Then most likely its a going to be a bad customer.

(Refer Slide Time: 33:09)

So, there are some prediction where it has said it is going to be a bad customer ok. Now I am going to do a table I call accuracy actual equals to data test dollar Good Bad and prediction predict equal to data test dollar Good Bad predicted values and this is the case. So, there are out of about 300, 212 cases where 202 plus 210 there are too many like bad customer is being missed and said that ok these are Good.

So, it is looks like the model is predicting most of the cases are as good 2000 out of 300 test cases. So, in the test cases test data set data test if I just do a simple table now data test comma good bad they are about 210 cases are good and 90 cases are bad. Now on the other hand if I just do a table of data test Good Bad pred. So, 282 cases are saying that they are Good and 18 cases they are saying Bad. So, this model is actually failing to you know failing to say who is Good who is Bad.

So, one thing people do try to understand that if I in case of 0.5 if I increase it to say 0.3 and will it be helpful. So, if you do that say probably yes, instead of 0.5 if you just reduce it now you have just too many probably saying or maybe 0.4. So, you have these many cases.

(Refer Slide Time: 36:07)

(Refer Slide Time: 36:30)



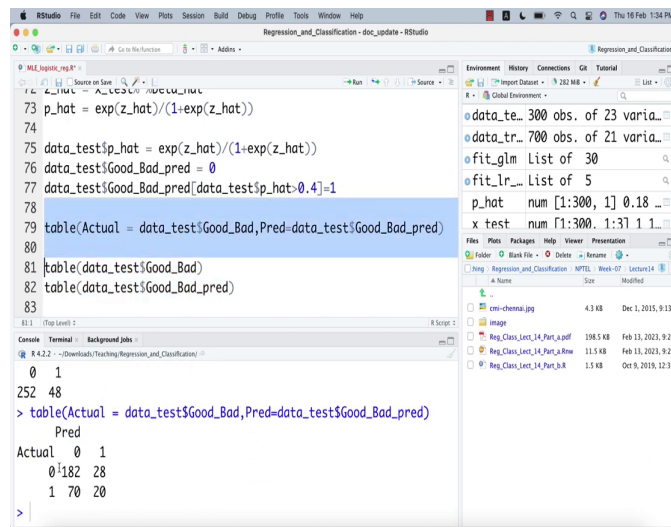So, if you have that now you have kind of 90 cases in actual cases these are the actual. So, 90 cases good and now it is better slightly better you can try out with this cut off marks. But the problem is another level now problem is there are two type of error one is truly there is a Good cases, which is being predicted as bad and there are Bad cases actually bad cases which is predicted as Good.

So, this is like actual bad cases which is predicted as good, which is actually bad thing you do not want this to happen if you as a bank manager you are giving loan to somebody and he or she is not returning in the money that is a really bad thing that is a real loss. This error on the other hand for a bank where actually the person is good, but bank thought ok if they will be bad customer and you are not giving loan.

So, in that case it is going to be a real problem. So, there are we are seeing in binary classification there are two kinds of problems. So, and one problem is typically called false positive another called false negative in statistical language it is called type 1 error and type 2 error. So, we are going to talk about that which one is good and which one is bad in the next lecture we are going to talk about that so, for now.

Thank you very much.