

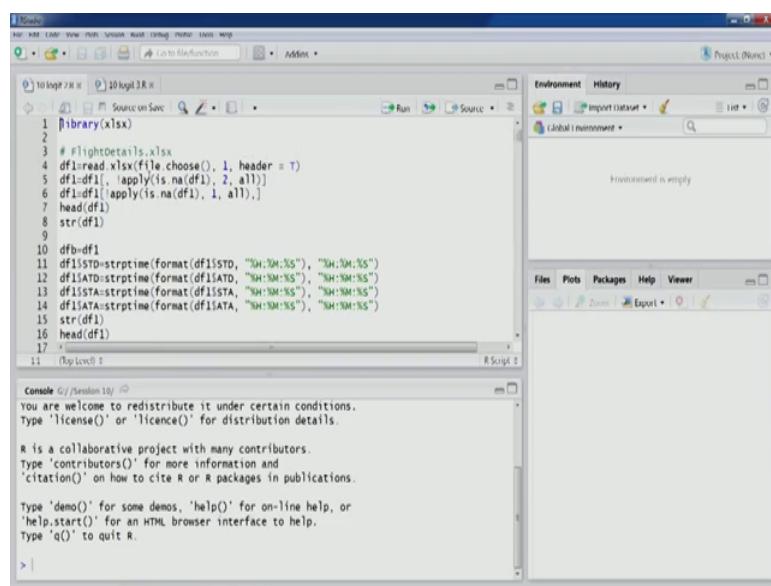
**Business Analytics & Data Mining Modeling Using R**  
**Dr. Gaurav Dixit**  
**Department of Management Studies**  
**Indian Institute of Technology, Roorkee**

**Lecture - 52**  
**Logistic Regression-Part VII**

Welcome to the course Business Analytics and Data Mining Modelling Using R. So, in previous few lectures we have been discussing different aspects of logistic regression so we will continue that discussion in this particular lecture as well.

So, in previous lecture we will be used apply details data set and then promotional offers data set as well. Some of the details regarding modelling exercise we could not cover.

(Refer Slide Time: 00:44)



The screenshot shows the RStudio interface with the following R code in the script editor:

```
library(xlsx)
# # flightDetails.xlsx
df1=read.xlsx(file.choose(), 1, header = T)
df1=df1[,apply(is.na(df1), 2, all)]
df1=df1[apply(is.na(df1), 1, all)]
head(df1)
str(df1)
dfb=df1
df1$STD=as.POSIXct(format(df1$STD, "%H.%M:%S"), "%H.%M:%S")
df1$ATA=as.POSIXct(format(df1$ATA, "%H.%M:%S"), "%H.%M:%S")
df1$ATA=as.POSIXct(format(df1$ATA, "%H.%M:%S"), "%H.%M:%S")
str(dfb)
head(dfb)
```

The console window displays the R startup message and the command prompt > |

So, we will do that and discuss. So, let us move to let us import this data set. So, we will let us import the library xlsx. So, we are again going to use this particular data set flight details.

(Refer Slide Time: 01:00)

The screenshot shows the RStudio interface. The top menu bar includes File, Edit, View, Plots, Source, Audit, Debug, Prefs, Help, and Help. The title bar says "RStudio". The main area has tabs for "10 logit.R" and "10 logit.R". Below the tabs is a toolbar with icons for Save, Run, Source, and others. The code editor contains the following R script:

```
library(xlsx)
# #FlightDetails.xlsx
df1<-read.xlsx(file.choose(), 1, header = T)
df1<-df1[, apply(is.na(df1), 2, all)]
df1<-df1[apply(is.na(df1), 1, all),]
head(df1)
str(df1)
dfb<-df1
df1$STD<-strptime(format(df1$STD, "%H:%M:%S"), "%H:%M:%S")
df2$STD<-strptime(format(df1$ATA, "%H:%M:%S"), "%H:%M:%S")
df3$STA<-strptime(format(df1$ATA, "%H:%M:%S"), "%H:%M:%S")
df4$ATA<-strptime(format(df1$ATA, "%H:%M:%S"), "%H:%M:%S")
str(dfb)
head(dfb)
#
```

The "Environment" tab on the right shows "Global environment" with an empty list. The "Files" tab is selected in the bottom navigation bar.

So, let us import this. Now, let us remove NA columns NA rows, let us look at these are the observations.

(Refer Slide Time: 01:18)

Let us the structure for the data frame. So, we will follow some of the steps that we have gone through in previous lectures as well, so we will just go through them.

(Refer Slide Time: 01:34)

```
23 df1=cbind(df1, DEPT)
24 df1$DEPT<-as.factor(df1$DEPT)
25 df1$Day<-as.factor(df1$Day)
26 levels(df1$Day)
27 levels(df1$Day)
28 levels(df1$Day)<-c("Sunday", "Monday")
29 df1$FLTIME<-as.difftime(as.character(df1$FLTIME))
30
31 str(df1)
32 head(df1)
33
34 dfb1=df1
35 dfb1[df1[,-c(1,3,5:8,10,12)]]
36 str(df1)
37 head(df1)
38
39
> dfb1=df1
> |
```

```
4 2017-08-28 08:43:00 2017-08-28 10:30:00 2017-08-28 10:24:00 BLR Monday
5 2017-08-28 02:25:00 2017-08-28 03:20:00 2017-08-28 04:06:00 BLR Monday
6 2017-08-28 02:52:00 2017-08-28 04:35:00 2017-08-28 04:32:00 BLR Monday
#> flight.status DIST FLTIME DEPT
1  ontime 842 84 mins 0-12
2  ontime 842 94 mins 0-12
3  delayed 842 79 mins 0-12
4  ontime 842 101 mins 0-12
5  delayed 842 101 mins 0-12
6  ontime 842 100 mins 0-12
> dfb1=df1
> |
```

Now, there is this particular exercise in the previous lecture we had used a separate grouping for departure time. So, again I have done certain changes into this grouping also, but this is not very important; however, let us run the model with a new grouping for departure time interval. So, this is the range we already familiar with.

(Refer Slide Time: 02:00)

```
19 range(df1$ATD)
20 breaks1=strptime("00:00:00", "%H:%M:%S")
21 breaks2=strptime("12:00:00", "%H:%M:%S")
22 DEPT=ifelse(df1$ATD>breaks2 & df1$ATD<=breaks1, "0-12", "12-24")
23
24 df1=cbind(df1, DEPT)
25 df1$DEPT<-as.factor(df1$DEPT)
26 df1$Day<-as.factor(df1$Day)
27 levels(df1$Day)
28 levels(df1$Day)<-c("Sunday", "Monday")
29 df1$FLTIME<-as.difftime(as.character(df1$FLTIME))
30
31 str(df1)
32 head(df1)
33
34 dfb1=df1
35
> dfb1=df1
> |
```

```
[1] "2017-08-28 00:40:00 IST" "2017-08-28 20:00:00 IST"
> breaks1=strptime("00:00:00", "%H:%M:%S")
> breaks2=strptime("12:00:00", "%H:%M:%S")
> DEPT=ifelse(df1$ATD>breaks2 & df1$ATD<=breaks1, "0-12", "12-24")
> df1=cbind(df1, DEPT)
> df1$DEPT<-as.factor(df1$DEPT)
> df1$Day<-as.factor(df1$Day)
> levels(df1$Day)
[1] "1"
> levels(df1$Day)<-c("Sunday", "Monday")
> df1$FLTIME<-as.difftime(as.character(df1$FLTIME))
> |
```

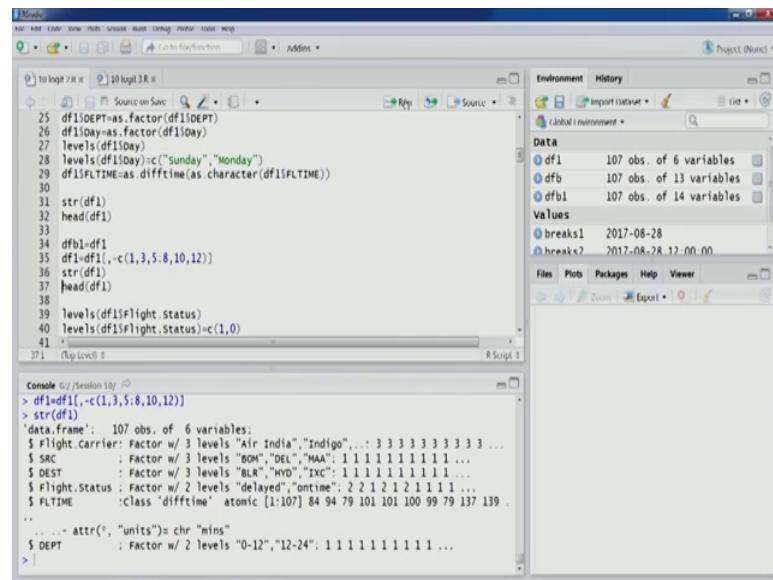
So, break is now 0 and 12, so these are the 2 breaks, 0 hours and 12 hours. Now this is how we are creating depart time variable. So, if less than a breaks 2, breaks 2, and breaks

1, then 0 to 12 so; that means, within if the timing is within these two you know hours 0 hours, and 12 hours.

Then first category that is 0 to 12 otherwise it is going to be the second category that is 12 to 14. So, let us create this variable let us append this to the data frame let us change it to a factor variable day variable as well let us cut the labels, change the labels flight time.

Let us also change it to appropriate format, now this is the structure that we have, now after taking backup we would not like to you know take forward some of the variables so let us get rid of them. Let us look at the structure again now these are the variables, now these are first few values for 6 values you can see everything is ok.

(Refer Slide Time: 03:04)



The screenshot shows the RStudio interface. The left pane contains an R script with the following code:

```
10 logit.R
25 df1$DEPT<-as.factor(df1$DEPT)
26 df1$Day<-as.factor(df1$Day)
27 levels(df1$Day)
28 levels(df1$Day)<-c("Sunday", "Monday")
29 df1$FLTTIME<-as.difftime(as.character(df1$FLTTIME))
30
31 str(df1)
32 head(df1)
33
34 df1<-df1
35 df1<-df1[,-c(1,3,5:8,10,12)]
36 str(df1)
37 head(df1)
38
39 levels(df1$flight status)
40 levels(df1$flight.status)<-c(1,0)
41
42 #skip Level 0
```

The right pane shows the Environment view with the following objects:

- Data:
  - df1: 107 obs. of 6 variables
  - dfb: 107 obs. of 13 variables
  - dfb1: 107 obs. of 14 variables
- Values:
  - breaks1: 2017-08-28
  - breaks2: 2017-08-28 12:00:00

Now, let us work on the outcome variable like we have been doing in previous lectures. So, let us change it to numeric code, let us change the reference category, now this is what it becomes now this is ok.

(Refer Slide Time: 03:30)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Shows R code for partitioning the dataset df1 into training (dftrain) and testing (dftest) partitions. It uses a 90%:10% split and prints the levels of the 'Flight.Status' factor.
- Console:** Displays the output of the R code, including the levels of the 'Flight.Status' factor for both training and testing datasets.
- Environment:** Shows the global environment with variables df1, dfb, dfb1, breaks1, and breaks2.

```
34 dfb1<-df1
35 dfb1<-df1[,-c(1,3,5:8,10:12)]
36 str(dfb1)
37 head(dfb1)
38
39 levels(df1$Flight.Status)
40 levels(df1$Flight.Status)<-(1,0)
41 head(df1$Flight.Status)
42 df1$Flight.Status<-relevel(df1$Flight.Status, ref = "0")
43 str(df1$Flight.Status)
44 head(df1$Flight.Status)
45
46 # Partitioning: 90%:10%
47 partidx<-sample(1:nrow(df1), 0.9*nrow(df1), replace = F)
48 dftrain<-df1[partidx,]
49 dftest<-df1[-partidx,]
50
51
52
53
54
55
```

```
[1] "delayed" "ontime"
> levels(df1$Flight.Status)<-(1,0)
> head(df1$Flight.Status)
[1] 0 0 1 0 1 0
Levels: 1 0
> df1$Flight.Status<-relevel(df1$Flight.Status, ref = "0")
> str(df1$Flight.Status)
Factor w/ 2 levels "0","1": 1 1 2 1 2 1 2 2 2 2 ...
> head(df1$Flight.Status)
[1] 0 0 1 0 1 0
Levels: 0 1
> |
```

Now, we can move ahead let us do our partitioning 2 partitions, training partition, testing partition and 90 percent for training partition and 10 of observation for test partition.

(Refer Slide Time: 03:40)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Shows R code for fitting a logistic regression model (mod3) using the glm function on the dftrain dataset.
- Console:** Displays the output of the R code, including the levels of the 'Flight.Status' factor for both training and testing datasets.
- Environment:** Shows the global environment with variables df1, dfb1, dftrain, dfbtest, and dfbtrain.

```
39 levels(df1$Flight.Status)
40 levels(df1$Flight.Status)<-(1,0)
41 head(df1$Flight.Status)
42 df1$Flight.Status<-relevel(df1$Flight.Status, ref = "0")
43 str(df1$Flight.Status)
44 head(df1$Flight.Status)
45
46 # Partitioning: 90%:10%
47 partidx<-sample(1:nrow(df1), 0.9*nrow(df1), replace = F)
48 dftrain<-df1[partidx,]
49 dftest<-df1[-partidx,]
50
51 mod3<-glm(Flight.Status ~ ., family = binomial(link = "logit"), data = dftrain)
52
53 #options(scipen=999)
54 summary(mod3)
55
```

```
[1] 0 0 1 0 1 0
Levels: 1 0
> df1$Flight.Status<-relevel(df1$Flight.Status, ref = "0")
> str(df1$Flight.Status)
Factor w/ 2 levels "0","1": 1 1 2 1 2 1 2 2 2 2 ...
> head(df1$Flight.Status)
[1] 0 0 1 0 1 0
Levels: 0 1
> partidx<-sample(1:nrow(df1), 0.9*nrow(df1), replace = F)
> dftrain<-df1[partidx,]
> dftest<-df1[-partidx,]
> |
```

Now, the same function glm that can be used to again model this. So, these are the results.

(Refer Slide Time: 03:53)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Contains R code for a logistic regression model. The code includes data partitioning, model fitting, and summary statistics.
- Console:** Displays the output of the R code, specifically the summary of the logistic regression model. The output table shows coefficients, standard errors, z-values, and p-values for various predictors.
- Environment:** Shows the global environment with objects like `dfb`, `dfb1`, `breaks1`, `breaks2`, `DEPT`, and `mod3`.

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-0.843187	0.861801	-0.978	0.37288
Flight.CarrierIndigo	-1.957116	0.747139	-2.619	0.00881 **
Flight.CarrierJet Airways	0.094405	0.598515	0.158	0.87467
SRCDEL	-0.121514	0.615738	-0.197	0.84356
SRCHAA	-0.877616	0.708888	-1.238	0.21571
DESTHYD	-1.024817	0.617856	-1.659	0.09718 .
DESTIXC	-0.845508	0.673242	-1.256	0.20916
FLTTIME	0.014289	0.004909	2.911	0.00361 **
DEPT12-24	1.077656	1.254554	0.859	0.39034
---				
Signif. codes:	0 '***'	0.001 '**'	0.01 '*'	0.05 '.'

Let us look at once again now you would see every time like we have been running this particular model on the same data set and every time. The significance levels have been changing that as I have been explaining smaller data set and therefore, subject to change in terms of as we as the observation that are part of training data set training partition this results will also slightly change and mainly with respect to significance level.

Now, you can see again flight carrier indigo has become significant to star level right and we can also see that destination has also become significant right. And then we can also see flight time has also you know modes you know higher level of significance.

However, more important thing is look at the p varies we can see that this one source for madras as well this is also smaller p values this is anyway significant flight carrier and destination anyway this is significant at ninety percent this is also significant. However, the new grouping that we have created out of department departure time intervals that also not comes out to be significant; however, p value is now smaller.

(Refer Slide Time: 05:14)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Displays R code for logistic regression. The code includes data partitioning, model fitting using `glm`, and summary statistics.
- Console:** Shows the output of the R code, including coefficient estimates for various flight carriers and other variables, and a dispersion parameter for the binomial family.
- Environment:** Shows the global environment with objects like `dfb`, `dfb1`, `breaks1`, `breaks2`, `DEPT`, and `mod3`.

So, with this we will discuss the important aspect of logistic regression so that is the majors of goodness of fit.

(Refer Slide Time: 05:19)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Displays R code for logistic regression, similar to the previous slide.
- Console:** Shows the output of the R code.
- Help:** A tooltip for the `glm` function is displayed, providing information about Generalized Linear Models.

So, just like the linear regression multiple linear regression logistic regression also is a statistical technique primarily. And therefore, in a statistical modelling the main objective is to fit to data as we have talked about this aspect many times before. So, therefore, in multiple linear regression we have a particular metric called a multiple R square and adjusted R square which are used to assess how good the model is fitting to the data.

So, similarly since logistic regression is also a statistical technique so there also we are required to have matrix for good to measure the goodness of it how well model is fitting to the data. So, what are those matrix so because there is certain key differences in logistic regression and linear regression. So, we will talk about some of those matrix now.

So, you can see in the code that I have created a vector here gf first one is mod 3 the model that we have just computed and you can see degree of residual degree of freedom. So, this df dot residual is one of the value that is returned by and the glm function and gives us the residual degree of freedom then we have deviance.

So, this is again the returns the deviance value returned by the glm function right then few other things which are mainly for the descriptive purposes. For example, this table result which is for the outcome variable here and then divided by the full number observation that will give us percentage success in training data and then we have iterations.

So, as we talked about the particular estimation technique that is used in that is used in logistic regression is different from multiple linear regression we can look at for more details we can look at here. So, we talked about that Emily is used for typically for used for logistic regression; however, you can see that for example, we have been using glm function, and within that if we go we look up for the some of the arguments.

(Refer Slide Time: 07:48)

```

50
51 mod3<-glm(flight.Status ~ ., family = binomial(link = "logit"), data = dftrain)
52
53 #options(scipen=999)
54 summary(mod3)
55
56 # Measures of Goodness of fit
57 gf<-c(mod3$df.residual, mod3$deviance,
58        100*table(dftrain$flight.Status)[[1]]/
59        length(dftrain$flight.Status),mod3$iter,
60        1-(mod3$deviance/mod3$null.deviance))
61 gf<-data.frame(gf, optional=TRUE)
62 rownames(gf)<-c("Residual df", "Std. Dev. Estimate",
63                  "% Success in training data", "#Iterations used",
64                  "Multiple R-squared")
65 gf
66

```

```

5942 Top Level: 0

```

```

Flight.CarrierJet Airways 0.094405 0.598515 0.158 0.87467
SRCDL -0.121514 0.615738 -0.197 0.84356
SRCHAA -0.877616 0.708888 -1.238 0.21571
DESTHYD -1.024817 0.617856 -1.659 0.09718
DESTXLC -0.845506 0.673242 -1.256 0.20916
FLTIME 0.014289 0.004909 2.911 0.00361 ==
DEPT12-24 1.077656 1.254554 0.859 0.39034
```
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
(Dispersion parameter for binomial family taken to be 1)

```

Environment pane:

- dfb: 107 obs. of 13 variables
- dfb1: 107 obs. of 14 variables
- Values:
  - breaks1: 2017-08-28
  - breaks2: 2017-08-28 12:00:00
  - DEPT: chr [1:107] "0-12" "0-12" ...
  - mod3: List of 30

Console pane:

```

glm(formula, family = gaussian, data, weights,
    na.action, start, nstart, min,
    control = list(...), model = TRUE, inc =
    FALSE, y = TRUE, contrasts = NULL)

glm.fit(y, w, weights = rep(1, nobs),
        start = NULL, cluster = NULL, mu =
        offset = rep(0, nobs), family = q,
        control = list(), intercept = TRUE)

## S3 method for class 'glm'
weights(object, type = c("prior", "working"))

```

Arguments:

formula: an object of class "formula" (or one

Specifically for this purpose the estimation technique purpose will get more detail. So, we will see that glm for dot fit.

(Refer Slide Time: 08:00)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Contains R code for fitting a logistic regression model. The code includes setting options, loading data, fitting the model, and printing the summary. It also calculates various measures of goodness of fit like AIC, BIC, and deviance.
- Console:** Displays the output of the R code, including the model summary, coefficient table, and dispersion parameter.
- Environment:** Shows the global environment with objects like `dfb`, `dfb1`, `breaks1`, `breaks2`, `DEPT`, and `mod3`.
- Help:** A tooltip for the `glm` function is open, explaining it as a function that uses iteratively reweighted least squares (IWLS) to fit generalized linear models.

We will see that glm dot fit method. So, this we can see iteratively re weighted least square is used the particular function that we are using glm iteratively re weighted least square function is used which is quite similar in approach with respect to Emily estimation technique that we talked about.

(Refer Slide Time: 08:41)

The slide has the following content:

## LOGISTIC REGRESSION

- Linear Regression for a categorical outcome variable?
  - Can be done by treating the outcome variable as continuous and coding it numerically
  - However, anomalies will lead to spurious modeling
    - Predictions can take any value, not just dummy values {0,1}
    - Outcome variable or residuals don't follow normal distribution
      - binomial distribution
    - Variance of outcome variable is not constant across all records (violation of homoscedasticity)
      - $np(1-p)$

So, MLM techniques, sorry MLM maximum likelihood method that we talked about in our discussion as we can look in these slides as well maximum likelihood method MLM method that we talked about. So, this is quite similar to what this the discussion that we have.

So, in particular R's implementation glm function that we are using it is the iteratively re weighted least square that estimation technique that is used to estimate the coefficients that we have been doing quite similar to a MLM and as we talked about that number of iterations have to be performed to reach to estimate these parameters so that we can get the best model which is the model best model which is fitting the data.

So, number of iterations actually indicate that and then we have one matrix which is quite similar to what we have multiple R square and linear regression right. So, this is actually computed where using the deviance value. So, null deviance and the standard deviation estimate or deviance that we have can look at the returned values here.

So, you can see deviance is the one of the return value so this is one and then we also null deviance is also a return. So, this is null deviance is come with respect to the naïve rule. So, 1 minus this deviance divided by null deviance gives us a value which is quite similar to what we have in multiple linear regression multiple R square.

So, this particular value will be will give us a metric which can be used to understand the goodness of fit follows logistic regression model and as on its own deviance also can be used it is quite similar to what we have there in as I see sum of squares error.

So, this is quite similar deviance is quite similar to that and then we can have one metric as I talked about similar to multiple R square. So, these metrics can be used to assess; the assess the fitness or model goodness of fitness goodness of fit of logistic regression model.

So, let us compute some of these things so, residual degree of freedom deviance which is similar to SSC in linear regression, then we have this proportions percentage in successive training data. Then we have number of iteration and then we have a multiple R square kind of metric let us compute this.

(Refer Slide Time: 11:34)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Displays R code for summarizing a model object named "mod". The code includes calculating residual degrees of freedom (df), standard deviation estimates, and other statistics.
- Console:** Shows the execution of the R code, resulting in a table of coefficients. The table includes columns for Estimate, Std. Error, z value, and Pr(>|z|). The first few rows of the table are:

|                           | Estimate  | Std. Error | z value | Pr(> z )   |
|---------------------------|-----------|------------|---------|------------|
| (Intercept)               | -0.843187 | 0.861801   | -0.978  | 0.32788    |
| Flight.CarrierIndigo      | -1.957116 | 0.747139   | -2.619  | 0.00881 ** |
| Flight.CarrierJet Airways | 0.094405  | 0.598515   | 0.158   | 0.87467    |

- Environment:** Shows the global environment with objects like "df", "dftrain", "dfb", and "dfb1". It also displays help documentation for "null.deviance" and "iter".

Let us create a data frame and row names we have given some and these are the values. So, you can see residual df df is 87 here, so as you can see let us again have a look df 1 training partition we have 96 observations.

(Refer Slide Time: 11:56)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Displays R code for summarizing a model object named "mod". The code includes calculating residual degrees of freedom (df), standard deviation estimates, and other statistics.
- Console:** Shows the execution of the R code, resulting in a table of coefficients. The table includes columns for Estimate, Std. Error, z value, and Pr(>|z|). The first few rows of the table are:

|                           | Estimate  | Std. Error | z value | Pr(> z )   |
|---------------------------|-----------|------------|---------|------------|
| (Intercept)               | -0.843187 | 0.861801   | -0.978  | 0.32788    |
| Flight.CarrierIndigo      | -1.957116 | 0.747139   | -2.619  | 0.00881 ** |
| Flight.CarrierJet Airways | 0.094405  | 0.598515   | 0.158   | 0.87467    |

- Environment:** Shows the global environment with objects like "df", "dftrain", "dfb", and "dfb1". It also displays help documentation for "null.deviance" and "iter".

And if we go back to our summary results right if we go back to our summary results we can see that how many variables we have here 1, 2, 3, 4, 5, 6, 7, 8; 8 variables. And we can see that 87 is the residual.

(Refer Slide Time: 12:06)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Contains R code for a logistic regression model. The code includes calculations for degrees of freedom, deviance, and R-squared.
- Console:** Displays the output of the R code, including the final R-squared value of 0.2044732.
- Environment:** Shows the global environment with various objects like df1, dftrain, dfb, dfb1, and gf.
- Help:** A tooltip for 'null.deviance' is open, explaining it as the deviance for the null model comparable with deviance. It notes that the null model will include the offset, and an intercept if there is one in the model. It also mentions that this will be incorrect if the link function depends on the data other than through the fitted mean; specifying a zero offset to force a correct calculation.

So, 87 plus 7 that makes it 94 that is  $n - 1$ , so that is the computation that is how the degrees of freedom have been computed, so this is a correct value here. And then we have deviance value which is also called as standard deviation estimated by some software's some statistical commercial statistical software's.

And then so this is all similar to what we have SSE sum of squares error in multiple linear regression then we have number of iteration that have been used to arrive at the particular model and that we have that is not 3 in this case. Then we have a value similar to multiple are squares we can see 20 percent, of the variability in the outcome variable has been explained by this model. So, all we talk about that this is being computed by  $1 - \frac{\text{mod3deviance}}{\text{null.deviance}}$ , divided by null deviance.

Now, whether so in terms of on further in terms of deviance the null deviance represents the naive rule value. So, we have to see how much our model has been able to how much reduction in deviance has been done by our model and whether that that is significant or not so that can also be that can also be performed using this chi square test that we can do.

So, we have one function `p.chisq`. So, there we can actually use these two values or we can take a difference of null deviance and deviance so that would be the reduction in deviance from a naive rule and you know how much deduction that our model has done.

And we can look at the number of predictors as degrees of freedom I could be used the number predictors that we use have used could be used degrees of freedom because these are freedom degree of freedom that had been used to reduce the deviance as we talked about from 95 minus 1 available degrees of freedom two we have reduced up to 87 that is residual degree of freedom so 7 predictors have been used.

So, that information can be used to perform chi square test and to find out the significance of whether the reduction has been significant or not. So, the third argument is lower tail that is specified as false and we can compute this chi square value so we can see that this is a small value. So, therefore, it seems that this reduction is significant which are also clear by the difference between the deviance and null deviance values.

So, we can also compute that we can see this is null deviance; this is null deviance. So, let us look at the value and we can have the deviance. Let us look at the value so you can see.

(Refer Slide Time: 15:17)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Displays R code for a logistic regression model. The code includes:
  - Reading data from 'flights' and creating a subset 'dftrain'.
  - Creating a formula 'mod3'.
  - Fitting the model with 'glm(mod3, family = binomial)'.
  - Checking the deviance with 'mod3\$deviance'.
  - Comparing the deviance with a null deviance using 'pchisq(mod3\$null.deviance, mod3\$deviance, length(mod3\$coefficients) - 1, lower.tail = F)'.
  - Outputting the result to the console.
- Console:** Shows the output of the R code, including the deviance values and the p-value from the chi-square test.
- Environment:** Shows the global environment with various objects and their characteristics.
- Help:** A tooltip for 'null.deviance' is displayed, explaining its meaning in the context of a GLM fit.

So, the difference is so there is good enough Reduction and deviance and that is why it also came as significant you know difference. So, these are some of the matrix that can actually be used to understand to a measure the goodness of fit to assess goodness of fit for a model as we talked about I talked about that in a statistical setting. So, these are some of the matrix which would be more useful.

So, in a statistical modelling we stopped at when we build the model using the training partition. So, typically all the observations are used that are present and then we look assess the model with respect to some of these some of these matrix. Now, let us move forward to our next discussion point in logistic regression so that is let us move forward so that is this particular point whether linear regression can be used for a categorical outcome variable right.

So, there are there are some situations where linear regression can be used as a categorical outcome variable which we will discuss later; however, right now we are discussing some of the more important points with respect to overall general applicability of linear regression for a categorical outcome variable.

So, can be done technically it can be applied so we can treat the outcome variable as continuous. So, the categorical outcome variable can be treated as continuous variable. So, we can essentially do the numeric coding and keep it as a numeric variable and technically we can apply we will get the results; however, the results would be meaningful or not that we need to understand.

So, technically it can be applied we can read the categorical outcome variable as continuous variable we can code it numerically so that can be done. However, there are going to be anomalies that would lead to spurious modelling so what could be some of these things.

So, number one predict predictions can take any value not just any values so for example, that binary logistic regression model that we have been performing for on some of the data sets so they are the our outcome variable it is it typically has two classes class 1, and class 0.

And so, the values the remaining variable we will take is to 0 for class 1 and 0 and 1 for class 1; however, when we apply a linear regression model to a categorical outcome variable the prediction can take any values any real value can be taken and not just the dummy values 0 and 1 so that is one challenge.

How do we map some of these predators values which can which can be any real value to the actual values of the outcome variable 0 and 1. Now, outcome variable or residuals do not follow normal distribution. So, as we have discussed during a linear regression

that this is one of the important assumption that dependent variable that is outcome variable all residuals should follow normal distribution, but that is not the case as we can understand that categorical outcome variable will have just 2 values 0 and 1.

So, it is actually it actually follows a binomial distribution so this particular assumption would also be violated; however, we talked about that because we are in data mining modelling context. So, where as we talked about even if for you know prediction purposes even if this particular you know assumption is violated in terms of prediction it might not much matter much because generally check performance on validation partition and test partitions; however, this case is different.

The deviation from normal distribution is much higher it is actually different distribution binomial distribution so that is one problem. Now, the another assumptions that we talked about in multiple linear regression is homoscedasticity; however, if we apply linear regression to an a to a categorical outcome variable this particular assumption would also be violated. The variance of outcome variable that we that we expect to be constant across all the records that is the homoscedasticity property so we for to apply multiple interrogation we want our outcome variable to follow this to have this property.

So, variance should be constant across all time; however, if we look at the variance for our categorical outcome variable it is going to be this particular value  $n$  times  $p$  into 1 minus  $p$  and as you can see because this is dependent on the value of  $p$ . So, therefore, the variance will change as the value of  $p$  changes.

So, when the value of a probability value is actually close to 0, then the variance would be on the as lower side and when the value of  $p$  is approaches 1 then the variance would be on the higher side. So, therefore, the variance will be will not be constant and it will be it will increase as the probability value you know in case is from 0 to 1.

So, so some of these are some of the problems that we can see in that we can directly understand and why a linear regression and cannot be applied to category outcome variable in a general sense and the problems that we can see here. So, what we will do we will do an exercise in R to understand the same thing to understand this particular aspect.

So, what we will do we will apply a linear regression model on a logistic partition and see the see its applicability and see some of the anomalies or violations that are that could be there. So, for this purpose we are going to use as you can see multiple here the comment is multiple linear regression model for a categorical response. So, promotion offers is the data set that we are going to use for this particular exercise. So, let us import the data set let us edit load.

(Refer Slide Time: 22:16)

The screenshot shows the RStudio interface. The left pane displays an R script titled '10 logit.R' with the following code:

```

130 segments(1,0,nrow(df1lift),df1lift$cumActualClass2[nrow(df1lift)],lty = 3)
131 legend(5, 2, inset = 0.005,
132         c("Cumulative Flight Status when sorted using predicted values",
133           "Cumulative Flight Status using average"),
134         lty = c(1,2), bty = "n", cex = 0.7, x.intersp = 0.3, y.intersp = 0.5)
135
136 # Multiple Linear Regression Model for a categorical response
137 # promoffers.xlsx
138 df2<-read.xlsx(file.choose(), 1, header = T)
139 df2<-df2[,capply(is.na(df2), 2, all1)]
140 df2<-df2[capply(is.na(df2), 2, all1),]
141 str(df2)
142
143 dfb2<-df2
144 dfb2<-df2[,c(1,3,7,9)]
145 str(dfb2)
146
147 (Top Level) 

```

The right pane shows the 'Environment' tab with objects like df1, df1test, df1train, dfb, dfb1, and gf listed. The bottom pane shows the 'Console' tab with output from the R session, including:

```

Std. Dev. Estimate      104.2420821
% Success in training data 57.2916667
#Iterations used          5.0000000
Multiple R-squared          0.2044732
> pchisq(mod3$null.deviance-mod3$deviance, length(mod3$coefficients)-1, lower.tail = F)
[1] 0.000767512
> mod3$null.deviance
[1] 131.0353
> mod3$deviance
[1] 104.2421
> df2<-read.xlsx(file.choose(), 1, header = T)
| 

```

So, once the observations have been loaded into environment we as we will see in the environment section we will do some of these steps I think it is has been loaded yes. So, df 2 we can see 5000 observation, 9 variables so let us remove NA columns, or NA rows if there are any let us look at the structure so this is the data set. So, we are already familiar with this.

(Refer Slide Time: 22:41)

```

131 legend(5, 2, inset = 0.005,
132   c("Cumulative Flight Status when sorted using predicted values",
133     "Cumulative Flight Status using average"),
134   lty = c(1,2), bty = "n", cex = 0.7, x.intersp = 0.3, y.intersp = 0.5)
135
136 # Multiple Linear Regression Model for a categorical response
137 # Promoffers.xlsx
138 df2<-read.xlsx(file.choose(), 1, header = T)
139 df2<-df2[, !apply(is.na(df2), 2, all)]
140 df2<-df2[!apply(is.na(df2), 2, all),]
141 str(df2)
142
143 df2<-df2[
144 df2<-df2[, c(1,3,7,9)]
145 str(df2)
146 #df2$Promoffer<-as.factor(df2$Promoffer)
147
148 # (Top Level) 5
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
709
710
711
712
713
714
715
716
717
718
719
719
720
721
722
723
724
725
726
727
728
729
729
730
731
732
733
734
735
736
737
738
739
739
740
741
742
743
744
745
746
747
748
749
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
769
770
771
772
773
774
775
776
777
778
779
779
780
781
782
783
784
785
786
787
787
788
789
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
848
849
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
878
879
879
880
881
882
883
884
885
886
887
887
888
889
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
909
910
911
912
913
914
915
916
917
917
918
919
920
921
922
923
924
925
926
927
928
929
929
930
931
932
933
934
935
936
937
938
939
939
940
941
942
943
944
945
946
947
947
948
949
950
951
952
953
954
955
956
957
958
959
959
960
961
962
963
964
965
966
967
968
969
969
970
971
972
973
974
975
976
977
978
979
979
980
981
982
983
984
985
986
987
987
988
989
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1009
1010
1011
1012
1013
1014
1015
1016
1017
1017
1018
1019
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1029
1030
1031
1032
1033
1034
1035
1036
1037
1037
1038
1039
1039
1040
1041
1042
1043
1044
1045
1046
1047
1047
1048
1049
1049
1050
1051
1052
1053
1054
1055
1056
1057
1057
1058
1059
1059
1060
1061
1062
1063
1064
1065
1066
1066
1067
1068
1068
1069
1070
1071
1072
1073
1073
1074
1075
1075
1076
1077
1077
1078
1079
1079
1080
1081
1082
1083
1084
1085
1085
1086
1087
1087
1088
1089
1089
1090
1091
1092
1093
1093
1094
1095
1095
1096
1097
1097
1098
1099
1099
1100
1101
1102
1103
1103
1104
1105
1105
1106
1107
1107
1108
1109
1109
1110
1111
1111
1112
1112
1113
1113
1114
1114
1115
1115
1116
1116
1117
1117
1118
1118
1119
1119
1120
1120
1121
1121
1122
1122
1123
1123
1124
1124
1125
1125
1126
1126
1127
1127
1128
1128
1129
1129
1130
1130
1131
1131
1132
1132
1133
1133
1134
1134
1135
1135
1136
1136
1137
1137
1138
1138
1139
1139
1140
1140
1141
1141
1142
1142
1143
1143
1144
1144
1145
1145
1146
1146
1147
1147
1148
1148
1149
1149
1150
1150
1151
1151
1152
1152
1153
1153
1154
1154
1155
1155
1156
1156
1157
1157
1158
1158
1159
1159
1160
1160
1161
1161
1162
1162
1163
1163
1164
1164
1165
1165
1166
1166
1167
1167
1168
1168
1169
1169
1170
1170
1171
1171
1172
1172
1173
1173
1174
1174
1175
1175
1176
1176
1177
1177
1178
1178
1179
1179
1180
1180
1181
1181
1182
1182
1183
1183
1184
1184
1185
1185
1186
1186
1187
1187
1188
1188
1189
1189
1190
1190
1191
1191
1192
1192
1193
1193
1194
1194
1195
1195
1196
1196
1197
1197
1198
1198
1199
1199
1200
1200
1201
1201
1202
1202
1203
1203
1204
1204
1205
1205
1206
1206
1207
1207
1208
1208
1209
1209
1210
1210
1211
1211
1212
1212
1213
1213
1214
1214
1215
1215
1216
1216
1217
1217
1218
1218
1219
1219
1220
1220
1221
1221
1222
1222
1223
1223
1224
1224
1225
1225
1226
1226
1227
1227
1228
1228
1229
1229
1230
1230
1231
1231
1232
1232
1233
1233
1234
1234
1235
1235
1236
1236
1237
1237
1238
1238
1239
1239
1240
1240
1241
1241
1242
1242
1243
1243
1244
1244
1245
1245
1246
1246
1247
1247
1248
1248
1249
1249
1250
1250
1251
1251
1252
1252
1253
1253
1254
1254
1255
1255
1256
1256
1257
1257
1258
1258
1259
1259
1260
1260
1261
1261
1262
1262
1263
1263
1264
1264
1265
1265
1266
1266
1267
1267
1268
1268
1269
1269
1270
1270
1271
1271
1272
1272
1273
1273
1274
1274
1275
1275
1276
1276
1277
1277
1278
1278
1279
1279
1280
1280
1281
1281
1282
1282
1283
1283
1284
1284
1285
1285
1286
1286
1287
1287
1288
1288
1289
1289
1290
1290
1291
1291
1292
1292
1293
1293
1294
1294
1295
1295
1296
1296
1297
1297
1298
1298
1299
1299
1300
1300
1301
1301
1302
1302
1303
1303
1304
1304
1305
1305
1306
1306
1307
1307
1308
1308
1309
1309
1310
1310
1311
1311
1312
1312
1313
1313
1314
1314
1315
1315
1316
1316
1317
1317
1318
1318
1319
1319
1320
1320
1321
1321
1322
1322
1323
1323
1324
1324
1325
1325
1326
1326
1327
1327
1328
1328
1329
1329
1330
1330
1331
1331
1332
1332
1333
1333
1334
1334
1335
1335
1336
1336
1337
1337
1338
1338
1339
1339
1340
1340
1341
1341
1342
1342
1343
1343
1344
1344
1345
1345
1346
1346
1347
1347
1348
1348
1349
1349
1350
1350
1351
1351
1352
1352
1353
1353
1354
1354
1355
1355
1356
1356
1357
1357
1358
1358
1359
1359
1360
1360
1361
1361
1362
1362
1363
1363
1364
1364
1365
1365
1366
1366
1367
1367
1368
1368
1369
1369
1370
1370
1371
1371
1372
1372
1373
1373
1374
1374
1375
1375
1376
1376
1377
1377
1378
1378
1379
1379
1380
1380
1381
1381
1382
1382
1383
1383
1384
1384
1385
1385
1386
1386
1387
1387
1388
1388
1389
1389
1390
1390
1391
1391
1392
1392
1393
1393
1394
1394
1395
1395
1396
1396
1397
1397
1398
1398
1399
1399
1400
1400
1401
1401
1402
1402
1403
1403
1404
1404
1405
1405
1406
1406
1407
1407
1408
1408
1409
1409
1410
1410
1411
1411
1412
1412
1413
1413
1414
1414
1415
1415
1416
1416
1417
1417
1418
1418
1419
1419
1420
1420
1421
1421
1422
1422
1423
1423
1424
1424
1425
1425
1426
1426
1427
1427
1428
1428
1429
1429
1430
1430
1431
1431
1432
1432
1433
1433
1434
1434
1435
1435
1436
1436
1437
1437
1438
1438
1439
1439
1440
1440
1441
1441
1442
1442
1443
1443
1444
1444
1445
1445
1446
1446
1447
1447
1448
1448
1449
1449
1450
1450
1451
1451
1452
1452
1453
1453
1454
1454
1455
1455
1456
1456
1457
1457
1458
1458
1459
1459
1460
1460
1461
1461
1462
1462
1463
1463
1464
1464
1465
1465
1466
1466
1467
1467
1468
1468
1469
1469
1470
1470
1471
1471
1472
1472
1473
1473
1474
1474
1475
1475
1476
1476
1477
1477
1478
1478
1479
1479
1480
1480
1481
1481
1482
1482
1483
1483
1484
1484
1485
1485
1486
1486
1487
1487
1488
1488
1489
1489
1490
1490
1491
1491
1492
1492
1493
1493
1494
1494
1495
1495
1496
1496
1497
1497
1498
1498
1499
1499
1500
1500
1501
1501
1502
1502
1503
1503
1504
1504
1505
1505
1506
1506
1507
1507
1508
1508
1509
1509
1510
1510
1511
1511
1512
1512
1513
1513
1514
1514
1515
1515
1516
1516
1517
1517
1518
1518
1519
1519
1520
1520
1521
1521
1522
1522
1523
1523
1524
1524
1525
1525
1526
1526
1527
1527
1528
1528
1529
1529
1530
1530
1531
1531
1532
1532
1533
1533
1534
1534
1535
1535
1536
1536
1537
1537
1538
1538
1539
1539
1540
1540
1541
1541
1542
1542
1543
1543
1544
1544
1545
1545
1546
1546
1547
1547
1548
1548
1549
1549
1550
1550
1551
1551
1552
1552
1553
1553
1554
1554
1555
1555
1556
1556
1557
1557
1558
1558
1559
1559
1560
1560
1561
1561
1562
1562
1563
1563
1564
1564
1565
1565
1566
1566
1567
1567
1568
1568
1569
1569
1570
1570
1571
1571
1572
1572
1573
1573
1574
1574
1575
1575
1576
1576
1577
1577
1578
1578
1579
1579
1580
1580
1581
1581
1582
1582
1583
1583
1584
1584
1585
1585
1586
1586
1587
1587
1588
1588
1589
1589
1590
1590
1591
1591
1592
1592
1593
1593
1594
1594
1595
1595
1596
1596
1597
1597
1598
1598
1599
1599
1600
1600
1601
1601
1602
1602
1603
1603
1604
1604
1605
1605
1606
1606
1607
1607
1608
1608
1609
1609
1610
1610
1611
1611
1612
1612
1613
1613
1614
1614
1615
1615
1616
1616
1617
1617
1618
1618
1619
1619
1620
1620
1621
1621
1622
1622
1623
1623
1624
1624
1625
1625
1626
1626
1627
1627
1628
1628
1629
1629
1630
1630
1631
1631
1632
1632
1633
1633
1634
1634
1635
1635
1636
1636
1637
1637
1638
1638
1639
1639
1640
1640
1641
1641
1642
1642
1643
1643
1644
1644
1645
1645
1646
1646
1647
1647
1648
1648
1649
1649
1650
1650
1651
1651
1652
1652
1653
1653
1654
1654
1655
1655
1656
1656
1657
1657
1658
1658
1659
1659
1660
1660
1661
1661
1662
1662
1663
1663
1664
1664
1665
1665
1666
1666
1667
1667
1668
1668
1669
1669
1670
1670
1671
1671
1672
1672
1673
1673
1674
1674
1675
1675
1676
1676
1677
1677
1678
1678
1679
1679
1680
1680
1681
1681
1682
1682
1683
1683
1684
1684
1685
1685
1686
1686
1687
1687
1688
1688
1689
1689
1690
1690
1691
1691
1692
1692
1693
1693
1694
1694
1695
1695
1696
1696
1697
1697
1698
1698
1699
1699
1700
1700
1701
1701
1702
1702
1703
1703
1704
1704
1705
1705
1706
1706
1707
1707
1708
1708
1709
1709
1710
1710
1711
1711
1712
1712
1713
1713
1714
1714
1715
1715
1716
1716
1717
1717
1718
1718
1719
1719
1720
1720
1721
1721
1722
1722
1723
1723
1724
1724
1725
1725
1726
1726
1727
1727
1728
1728
1729
1729
1730
1730
1731
1731
1732
1732
1733
1733
1734
1734
1735
1735
1736
1736
1737
1737
1738
1738
1739
1739
1740
1740
1741
1741
1742
1742
1743
1743
1744
1744
1745
1745
1746
1746
1747
1747
1748
1748
1749
1749
1750
1750
1751
1751
1752
1752
1753
1753
1754
1754
1755
1755
1756
1756
1757
1757
1758
1758
1759
1759
1760
1760
1761
1761
1762
1762
1763
1763
1764
1764
1765
1765
1766
1766
1
```

So, the variable that we are going to use for our outcome variable is going to be the promotional offer as you can see that we have commented out the lines of code, which we used earlier to convert these numeric variables going to convert numeric variables into the categorical variable. So, promotional offer an online so they are actually categorical variable factor variable, but we are not converting them into factor variable because we are going to apply a linear regression modelling. So, we will give them as numeric and we will apply.

So, a partitioning is the same 60 percent, 40 percent in this case so we can see that let us do partitioning. So, df 2 train you can see 3000 observations, 4 variables. Now the same align function is going to be used now the promotional offer is going to be request against all the predictors that are present in this particular dataset df 2 train. So, let us run this.

(Refer Slide Time: 24:30)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Contains R code for data partitioning, model fitting, and summary statistics. The code includes:
 

```

148 #str(df2)
149
150 # Partitioning: Tr:Te->60%:40%
151 partidx<-sample(1:nrow(df2), 0.6*nrow(df2), replace = F)
152 df2train<-df2[partidx,]
153
154 mod4<-lm(Promoffer ~ ., df2train)
155
156 #options(scipen=999)
157
158 mod4summ<-summary(mod4); mod4sum
159 mod4ava<-anova(mod4); mod4ava
160
161 DF4c<-mod4summ$statistic[["numDF"]],
162       mod4summ$statistic[["denDF"]],
163
164
165 #logLik(mod4)
      
```
- Console:** Displays the results of the R code execution, including:
 

|        | Min      | Q1       | Median   | Q3      | Max     |
|--------|----------|----------|----------|---------|---------|
| Income | -0.58172 | -0.13744 | -0.03657 | 0.05929 | 0.92044 |

Coefficients:

|             | Estimate   | Std. Error | t value | Pr(> t )   |
|-------------|------------|------------|---------|------------|
| (Intercept) | -0.2388135 | 0.0148069  | -16.129 | <2e-16 *** |
| Income      | 0.0033865  | 0.0001005  | 31.706  | <2e-16 *** |
| Family.size | 0.0372885  | 0.0040443  | 9.220   | <2e-16 *** |
| online      | -0.0058494 | 0.0093219  | -0.627  | 0.53       |

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' '
- Environment:** Shows global variables and their values.
- Help:** A tooltip for the `deviance` function is displayed.

Now, what we will do we will look at the summary table. Let us look at the results. So, as we can see that income is intercept significant income is significant and family size is significant we can look at the different estimates for example, income quite a small value from this family size is also 0.03, online is not significant it has not only found to significant as we can see.

(Refer Slide Time: 24:55)

The screenshot shows the RStudio interface. The top menu bar includes File, Edit, View, Plots, Session, Tools, Debug, Project, Help, and a 'Project (None)' dropdown. The main window has tabs for '10 logit 7.R' and '10 logit 3.R'. Below the tabs is a toolbar with icons for Source on Save, Run, Source, and Addins. The code editor contains R code for model selection and summary statistics. A message box is displayed in the center, indicating that the 'mod4sum4' object is a list of 11 elements. The right panel shows the Environment tab with variables like 'breaks1', 'breaks2', 'DEPT', 'mod3', 'mod4', and 'mod4sum4'. The bottom right corner displays a note about the deviance of the null model.

```
148 #str(df2)
149
150 # Partitioning: Tr:Te->60%:40%
151 partidx=sample(1:nrow(df2), 0.6*nrow(df2), replace = F)
152 df2train=df2[partidx,]
153
154 mod4=lm(prmoffer ~ ., df2train)
155
156 #options()
157 #options(scipen=999)
158
159 mod4sum4=summary(mod4); mod4sum4
160 mod4ava=anova(mod4); mod4ava
161
162 DF<-c(mod4sum4$fstatistic[["numdf"]], 
163 mod4sum4$fstatistic[["dendf"]])
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
```

(Top Level) R Script I

Console G:/Session 10>

```
estimate std. error t value Pr(>|t|)
(Intercept) -0.2388135 0.0148069 -16.129 <2e-16 ***
Income 0.0031865 0.0001005 33.706 <2e-16 ***
Family.size 0.0372885 0.0040443 9.220 <2e-16 ***
online -0.0058494 0.0093219 -0.627 0.53
...
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.251 on 2996 degrees of freedom
Multiple R-squared: 0.2774, Adjusted R-squared: 0.2767
F-statistic: 383.4 on 3 and 2996 DF, p-value: < 2.2e-16
```

Environment History

values

- breaks1 2017-08-28
- breaks2 2017-08-28 12:00:00
- DEPT chr [1:107] "0-12" "0-12" ...
- mod3 List of 30
- mod4 Large lm (12 elements, 80L)
- mod4sum4 List of 11

Files Plots Packages Help Viewer

H: fitting Generalized Linear Models | Find in Topic

families fitted by quasi-likelihood the value is NA.

null.deviance The deviance for the null model, comparable with deviance. The null model will include the offset, and an intercept if there is one in the model. Note that this will be incorrect if the link function depends on the data other than through the fitted mean specify a zero offset to force a correct calculation.

iter the number of iterations of NLVs used.

And we can have a look at the other values we can see adjusted R square R square and multiple R square. So, we can see if we you know 27 percent multiple R square value we look at the p well it is quite as small. So, the model is significant ah; however, as we talked about certain problems as we have discussed could be there certain anomalies could be there.

(Refer Slide Time: 25:21)

So, what we will do? We will look at computing some of those things to check whether those anomalies are there in this particular case. So, what we will do we are

running and over to extract some of the parameters here. So, we can see so sum of square errors mean of square errors and F value of a statistic and the probability values for these predictors income family size is given there in the ANOVA table.

So, what we are going to do is we are going to compute these particular values in a in a format that can be used for tabular presentation later on. So, mod 4 so on; we have a F statistics value. So, that is written as part of the summary function of models.

(Refer Slide Time: 26:19)

```

158 mod4sum<-summary(mod4); mod4sum
159 mod4ava<-anova(mod4); mod4ava
160
161 DF<-c(mod4$sumsq$statistic[["numdf"]],
162 mod4$sumsq$statistic[["dendf"]], 
163 mod4$sumsq$statistic[["numdf"]]+mod4$sumsq$statistic[["dendf"]])
164
165 SS<-c(sum(head(mod4ava[,"Sum Sq"],-1)),
166 mod4ava[,"Residuals","Sum Sq"],
167 sum(mod4ava[,"Sum Sq"]))
168
169 MS<-c(mean(head(mod4ava[,"Mean Sq"],-1)),
170 mod4ava[,"Residuals","Mean Sq"])
171
172 Fstat<-c(mod4$sumsq$statistic[["value"]])
173 Fstat
174
175 Fstat<-Fstat/mod4$sumsq$statistic[["value"]]
176
177 Fstat
178
179 Fstat
180
181 Fstat
182
183 Fstat
184
185 Fstat
186
187 Fstat
188
189 Fstat
190
191 Fstat
192
193 Fstat
194
195 Fstat
196
197 Fstat
198
199 Fstat
200
201 Fstat
202
203 Fstat
204
205 Fstat
206
207 Fstat
208
209 Fstat
210
211 Fstat
212
213 Fstat
214
215 Fstat
216
217 Fstat
218
219 Fstat
220
221 Fstat
222
223 Fstat
224
225 Fstat
226
227 Fstat
228
229 Fstat
230
231 Fstat
232
233 Fstat
234
235 Fstat
236
237 Fstat
238
239 Fstat
240
241 Fstat
242
243 Fstat
244
245 Fstat
246
247 Fstat
248
249 Fstat
250
251 Fstat
252
253 Fstat
254
255 Fstat
256
257 Fstat
258
259 Fstat
260
261 Fstat
262
263 Fstat
264
265 Fstat
266
267 Fstat
268
269 Fstat
270
271 Fstat
272
273 Fstat
274
275 Fstat
276
277 Fstat
278
279 Fstat
280
281 Fstat
282
283 Fstat
284
285 Fstat
286
287 Fstat
288
289 Fstat
290
291 Fstat
292
293 Fstat
294
295 Fstat
296
297 Fstat
298
299 Fstat
300
301 Fstat
302
303 Fstat
304
305 Fstat
306
307 Fstat
308
309 Fstat
310
311 Fstat
312
313 Fstat
314
315 Fstat
316
317 Fstat
318
319 Fstat
320
321 Fstat
322
323 Fstat
324
325 Fstat
326
327 Fstat
328
329 Fstat
330
331 Fstat
332
333 Fstat
334
335 Fstat
336
337 Fstat
338
339 Fstat
340
341 Fstat
342
343 Fstat
344
345 Fstat
346
347 Fstat
348
349 Fstat
350
351 Fstat
352
353 Fstat
354
355 Fstat
356
357 Fstat
358
359 Fstat
360
361 Fstat
362
363 Fstat
364
365 Fstat
366
367 Fstat
368
369 Fstat
370
371 Fstat
372
373 Fstat
374
375 Fstat
376
377 Fstat
378
379 Fstat
380
381 Fstat
382
383 Fstat
384
385 Fstat
386
387 Fstat
388
389 Fstat
390
391 Fstat
392
393 Fstat
394
395 Fstat
396
397 Fstat
398
399 Fstat
400
401 Fstat
402
403 Fstat
404
405 Fstat
406
407 Fstat
408
409 Fstat
410
411 Fstat
412
413 Fstat
414
415 Fstat
416
417 Fstat
418
419 Fstat
420
421 Fstat
422
423 Fstat
424
425 Fstat
426
427 Fstat
428
429 Fstat
430
431 Fstat
432
433 Fstat
434
435 Fstat
436
437 Fstat
438
439 Fstat
440
441 Fstat
442
443 Fstat
444
445 Fstat
446
447 Fstat
448
449 Fstat
450
451 Fstat
452
453 Fstat
454
455 Fstat
456
457 Fstat
458
459 Fstat
460
461 Fstat
462
463 Fstat
464
465 Fstat
466
467 Fstat
468
469 Fstat
470
471 Fstat
472
473 Fstat
474
475 Fstat
476
477 Fstat
478
479 Fstat
480
481 Fstat
482
483 Fstat
484
485 Fstat
486
487 Fstat
488
489 Fstat
490
491 Fstat
492
493 Fstat
494
495 Fstat
496
497 Fstat
498
499 Fstat
500
501 Fstat
502
503 Fstat
504
505 Fstat
506
507 Fstat
508
509 Fstat
510
511 Fstat
512
513 Fstat
514
515 Fstat
516
517 Fstat
518
519 Fstat
520
521 Fstat
522
523 Fstat
524
525 Fstat
526
527 Fstat
528
529 Fstat
530
531 Fstat
532
533 Fstat
534
535 Fstat
536
537 Fstat
538
539 Fstat
540
541 Fstat
542
543 Fstat
544
545 Fstat
546
547 Fstat
548
549 Fstat
550
551 Fstat
552
553 Fstat
554
555 Fstat
556
557 Fstat
558
559 Fstat
560
561 Fstat
562
563 Fstat
564
565 Fstat
566
567 Fstat
568
569 Fstat
570
571 Fstat
572
573 Fstat
574
575 Fstat
576
577 Fstat
578
579 Fstat
580
581 Fstat
582
583 Fstat
584
585 Fstat
586
587 Fstat
588
589 Fstat
590
591 Fstat
592
593 Fstat
594
595 Fstat
596
597 Fstat
598
599 Fstat
599>

```

So, when we apply summary on the model object so you can see this is nothing, but F statistic value for the model and then we have degree of freedom and here. So, we can see the degree of freedom residuals degree of freedom and the residuals degree of freedom and the regression degree of freedom here so that is going to be stored in this so in this data frame.

First we have the regression degree of freedom, then we have the residual degree of freedom and then we have the total so this data frame is about degree of freedom as is sustained by its named DF. Then we will compute some other square error so first we can see that in the ANOVA table, from the ANOVA table we are trying to extract out this these values.

(Refer Slide Time: 27:21)

The screenshot shows the RStudio interface. The left pane displays an R script with the following code:

```
162 DF<-c(mod4sum$statistic[["numdf"]],
163 mod4sum$statistic[["dendf"]],
164 mod4sum$statistic[["numdf"]]+mod4sum$statistic[["dendf"]])
165
166 SS<-c(sum(head(mod4ava[, "Sum Sq"], -1)),
167 mod4ava[["Residuals", "Sum Sq"]],
168 sum(mod4ava[, "Sum Sq"]))
169
170 MS<-c(mean(head(mod4ava[, "Mean Sq"], -1)),
171 mod4ava[["Residuals", "Mean Sq"], ""))
172
173 Fstatistic<-c(mod4sum$statistic[["value"]], "", "")
174 P<-pf(mod4sum$statistic[[1]], mod4sum$statistic[[2]], mod4sum$statistic[[3]])
175 lower.tail = F)
176 pvalue<-c(P, "", "")
177
178
179
180
```

The right pane shows the help page for the `F` distribution, specifically the `FDist` function. The description and usage sections are visible.

So, first some other square for the you know regressors regression and then for the residuals, and then total. So, this would be recorded in this particular variable SS. And then we look at the mean or mean of square errors, so that is also being extracted from the ANOVA table results, you can see this particular column mean square and these values are being instructed.

So, first 3 values so this head function as you can see now this role of this head function is quite different and both these computations we have used head function.

So, you can see the second argument is minus 1 so what it actually does is it gives us all the values except the last value in the vector. So, for example, mean square or some other square. So, there are 4 values and these 2 columns, and these 2 vectors.

So, except the last values that is corresponding to residuals the first 3 values are going to be written; that means, last value will be left out and the remaining and remaining n minus 1 values are going to be written. So, that is what we want. So, that will give us the corresponding sum of squares or mean sum of squares are for the deviation.

So, first that then residuals so let us compute this. Then let us also extract the F statistics, from the this particular vector that we have already seen. So, let us do that then what we are trying to do we are trying to compute the corresponding probability value.

So, probability value corresponding to the F for test so that that is how it is being committed pf is the function that can be used for more detail you can go into the help section and find out more information about pf. So, you can see this is F distribution.

(Refer Slide Time: 29:09)

The screenshot shows the RStudio interface. On the left, the 'Console' tab displays R code for calculating F statistics and p-values. The code involves summing residuals squared, calculating mean squares, and then using the mod4aval and mod4summfstatistic functions to compute the F statistic and its p-value. It also creates an anovaadf data frame with columns for Regression, Error, and Total. The 'Script' tab shows the same code. On the right, the 'Help' tab is open, displaying the documentation for the pf function. The documentation includes the function's name, a brief description of random generation for the F distribution, usage examples, and arguments. The arguments section defines 'q' as a vector of quantiles and 'p' as a vector of probabilities.

```

167 mod4aval("Residuals","Sum Sq"),
168 sum(mod4aval,"Sum Sq"))
169
170 MS<-c(mean(head(mod4aval[, "Mean Sq"], -1)),
171 mod4aval("Residuals","Mean Sq"), "")
172
173 Fstatistic<-c(mod4summfstatistic[("value")], "", "")
174 P<-pf(mod4summfstatistic[[1]], mod4summfstatistic[[2]], mod4summfstatistic[[3]],
175 lower.tail = F)
176 pvalue<-c(P, "", "")
177
178 anovaadf<-data.frame(DF, SS, MS, Fstatistic, pvalue)
179 rownames(anovaadf)<-c("Regression", "Error", "Total")
180 anovaadf
181
182 # Prediction for a new observation.
183
184
185
186
187
188
189
190
191 #> [1] 5

```

**Help on pf**

**Usage**

```

pf(r, df1, df2, ncp, log = FALSE)
pifq, dlf, dlf, ncp, lower.tail = TRUE, l
qf(p, df1, df2, ncp, lower.tail = TRUE, l
rlf, dlf, dlf, ncp)

```

**Arguments**

- q** vector of quantiles.
- p** vector of probabilities.

So, within F distribution we have F function which can be used to compute the corresponding p value. So, what we need is F statistics and the first argument then we need degree of freedom as second argument that is responding to a regression here and the second is then last one is residuals degree of freedom third argument and then we have a specified lower tail as false.

So, we will get the p value corresponding p value. So, let us record it in this format once this is done we can create this table data frame. And let us assign a names for this and let us have a look at this particular table. So, we can see that now once computed.

(Refer Slide Time: 29:50)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Contains R code for calculating ANOVA and F-statistics. The code includes calculations for degrees of freedom (DF), sum of squares (SS), mean square (MS), F statistic, and p-value.
- Console:** Displays the output of the R code, showing the results of the ANOVA table.
- Environment:** Shows the global environment with various objects defined, such as `anovaadf` (3 obs. of 5 variables), `df1` (107 obs. of 6 variables), etc.
- Help:** Shows the help documentation for the `pf` function, which is used for the F-distribution.

```
R Script Editor:
170 MS<-c(mean(head(mod4aval[, "Mean Sq"], -1)),  
171 mod4aval["Residuals","Mean Sq"], "")  
172  
173 Fstatistic<-c(mod4summsfstatistic[["value"]], "", "")  
174 P<-pf(mod4summsfstatistic[[1]], mod4summsfstatistic[[2]], mod4summsfstatistic[[3]])  
175 lower.tail = F  
176 pvalue<-c(P, "")  
177  
178 anovadf<-data.frame(DF, SS, MS, Fstatistic, pvalue)  
179 rownames(anovadf)<-c("Regression", "Error", "Total")  
180 anovadf  
181  
182 # Prediction for a new observation:  
183 # Annual income of 5 Lakhs with two family members  
184 # who is not active online  
185 predict(mod4, data.frame(income=5, Family.Size=2, Online=0))  
186  
187 (Copy Level) 8
```

```
Console [1/Session 10]:  
> Fstatistic<-c(mod4summsfstatistic[["value"]], "", "")  
> P<-pf(mod4summsfstatistic[[1]], mod4summsfstatistic[[2]], mod4summsfstatistic[[3]],  
+ lower.tail = F)  
+ pvalue<-c(P, "")  
> anovadf<-data.frame(DF, SS, MS, Fstatistic, pvalue)  
> rownames(anovadf)<-c("Regression", "Error", "Total")  
> anovadf
```

|            | DF   | SS        | MS                | Fstatistic       | pvalue                |
|------------|------|-----------|-------------------|------------------|-----------------------|
| Regression | 3    | 72.44317  | 24.1477243151055  | 383.361202942226 | 1.00014894386617e-210 |
| Error      | 2996 | 188.71649 | 0.062989483885631 |                  |                       |
| Total      | 2999 | 261.15967 |                   |                  |                       |

So, first column we have degrees of freedom. So, for regression there are 3 and then because as remember that we have used just 4 variable, so one being outcome variables, so 3 predictors. So, therefore, degree of freedom degree of freedom is here is 3 for regression then residuals it is the remaining that is n minus 1 is 2999, total number of observation and is 3000.

So, residual degree of freedom 2996 sum of a square for regression and for residual is also there. So, you can see that residual sum of square is much higher so that also is we can we saw that the variance was also on the lower side right earlier we saw that the variance that we had computed was on the lower side so that is also indicated here, mean square of errors is also there and then we have F statis F statistic and then we have p value is small value.

(Refer Slide Time: 31:00)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Contains R code for performing ANOVA and creating a histogram of residuals. The code includes setting up data frames, calculating F-statistics, and performing predictions.
- Console:** Displays the output of the R code, including the ANOVA table and a histogram of residuals.
- Environment:** Shows the global environment with objects like `anovaadf`, `df1`, `df2`, etc.
- Help:** Shows the help documentation for the `pf` function, which is used for random generation from the F-distribution.

```
178 anovadf<-data.frame(DF, SS, MS, Fstatistic, pvalue)
179 rownames(anovadf)<-c("Regression", "Error", "Total")
180 anovadf
181
182 # Prediction for a new observation:
183 # Annual income of 5 lakhs with two family members
184 # who is not active online
185 predict(mod4, data.frame(Income=5, Family.Size=2, Online=0))
186 # Set of values for Promoffer
187 #> df2<-data.table(df2$Promoffer)
188
189 range(mod4$residuals)
190 hist(mod4$residuals, main="", xlab="Residuals", xlim = c(-1,1))
191
192 #cleanup
193 rm(list = ls())
194
195 #> mod4<-NULL
```

```
Console | R Session 10/1
> Fstatistic<-c(mod4$sum$Fstatistic[["value"]], "", "")
> Pvalue<-c(mod4$sum$Fstatistic[[1]], mod4$sum$Fstatistic[[2]], mod4$sum$Fstatistic[[3]],
+ lower.tail = F)
> pvalue<-c(P, "", "")
> anovadf<-data.frame(DF, SS, MS, Fstatistic, pvalue)
> rownames(anovadf)<-c("Regression", "Error", "Total")
> anovadf
      DF     SS       MS   Fstatistic    pvalue
Regression  3 72.44317 24.1477243151055 383.361202942226 1.00014894386617e-210
Error      96 2996.188 71649.0 0.062989483885631
Total      99 2999.261.15967
```

So, this gives us some information about the model and now what we will do is we will use this model to the score of a particular observation. So, let us do a prediction for a new observation let us say this is our new observation is annual income of rupees 5 lakhs with 2 family members who is not active online.

So, this is information about the particular customers customer whether is going to accept the promotional offer or not. So, annual income is 5 lakh family size is 2 and not acting online. So, you can use the predict function first argument is as usual the model object mod 4 then in a data frame we are trying to pass on the values operate predictors for example, income 5.

So, it should be in the same unit as was used for the modelling exercise in the training partition. So, you can see 5 and family size is 2 and then online is because this particular customer is not active so 0. So, once we do this we will be able to predict this particular observation. So, you can see this comes out to be a negative value. Now that was one of the first point that we discussed in the slide.

(Refer Slide Time: 32:11)

## LOGISTIC REGRESSION

- Logistic Regression for Profiling Task
  - Apart from model performance on validation partition
  - Model's fit to data is assessed on training partition
    - However, still avoid overfitting
    - Usefulness of predictors is examined
  - Goodness of fit metrics
    - Overall fit of the model
      - Deviance (equivalent to SSE in linear regression)
      - 1 - Deviance/Null Deviance (equivalent to multiple R<sup>2</sup> in linear regression)
    - Single predictors

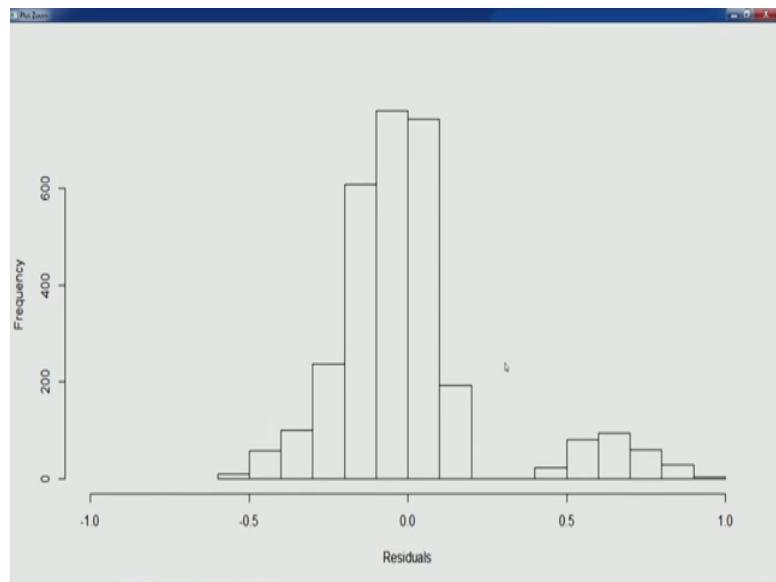
IIT Kharagpur | NPTEL ONLINE CERTIFICATION COURSE

15

Let us go back you can see here that the predictions can take any value not just dummy values. So, we can see that a negative value has been taken here and so that is in one anomaly that we can clearly see. And set of values for them was no offer that we already know 0 and 1, and the value is that comes out to be predicted value comes out to be minus 0.14 now.

So, one difficulty is how do we do our classification in this case? Now, let us look at the residuals so second anomaly that I discussed that outcome variable of residuals do not follow normal distribution. So, let us look at that let us plot a histogram and find out whether this is being followed.

(Refer Slide Time: 33:07)



You can see residuals when we plot residuals when we get histogram you can see this is clearly not being followed normal distribution is not being followed one grouping here and other grouping where this is of this lower values, lower frequency and this is higher frequency. So, this is clearly looks binomial or different groups kind of thing. So, this definitely not following so normal distribution and distortions due to real binomial distribution can be seen here.

So, the typically what exercises and the discussion that we have been doing was pertaining to classification tasks and as we talked about that this is a statistical technique logistic regression being a statistical technique is also applied in so used in statistical modelling. So, the kind of task that we generally do in a statistical modelling are quite similar to what we can call profiling tasks.

So, as we talked about in the you know starting lecture of logistic regression and that it is about understanding similarities and differences between two groups. So, logistic regression can also be used in can also use to understand what are the variables which you know which can bring out some of the similarities or differences between group, so let us discuss that aspect as well.

So, in profiling tasks when we talk about profiling tasks the situation is slightly different. So, in the classification task we typically build our model and look at the performance of that classifier that particular model using the classification matrix using the overall

accuracy or overall error matrix and you know some deviations when we have a class of interest.

So, those are the things that we typically do we also typically look at left chart especially when we have a class of interest to in left chart and Decile chart to see whether it is still despite you know higher error whether it still the model is useful in class of interest with respect to naive you know naive rule or a or a average case.

So, some of those things that we do in classification tasks; however, in profiling tasks what we do in classification we follow that. So, apart from model performance on validation partition we also asses models fit to data on training partition right because as I taught what in a statistical technique typically the whole sample is used; however, since we are using training partition for model building.

So, the models fit to data is assessed on training partition; however, model performance is assessed on validation partition for profiling tasks and models fit to data assessed on training partition. So, some of these things we talked about when we talk about goodness of fit measures we talked about the deviance, we talked about 1 minus deviance divided by null deviance that is equivalent of multiple R square some of those things you can see here as well.

So, models fit to data is assessed on training partition ah; however, we still focus on avoiding over fitting because as we talked about when we have matrix which look for and when we do modelling to achieve the you know goodness of fit then it can it can lead to over fitting. So, it still we would like to avoid over fitting and still be able to you know check the performance still be able to have good classification performance as well.

So, usefulness of usefulness of predictors is also examined in this particular case profiling so because it is about understanding similarities and differences between 2 groups. So, therefore, which predictor is more helpful in terms of bringing out those differences those similarities so when we build our model. We also look at the significance levels of some of these predictors we look at which predictors are significant which are not significant whether the not significant predictors can be dropped from the model.

And this has also should be looked at from the perspective of model performance because as we have taught about data mining model we would like to keep the insignificant variables also in the order if they provide some practical importance in terms of scoring new observation.

How and logistic statistical modelling we just drop the insignificant variables because we are just interested in understanding the phenomena, understanding the underlying relationship. So, profiling is quite similar to you know that that approach; however, because we are doing data mining modelling. So, we have to balance between these two, we have to balance between performance and also the main profiling tasks.

So, we have to really see whether the predictors can be dropped just like in statistical training or whether they have to be kept in the model because they also provide some practical significance for scoring new observation so that balance has to be achieved.

So, we have to avoid over fitting, we have to look for using usefulness or predictors in both the context data mining context prediction point of view and also statistical context in terms of understanding you know finding understanding variables which differentiate the groups, which bring out similarity or differences between groups.

So, this kind of exercise is done in profiling tasks and goodness of fit matrix that through an exercise in R that we have already understood overall fit. So, we look to understand first we look to understand the overall fit of the model, so if the overall fit of the model is good only then we go ahead and look at the individual variables.

So, first step typically is in profiling on a statistical modelling we look at the overall fit of the model. So, in this particular case logistic regression as we talked about the deviance is the metric we taught one previous lecture that could be used and we also said that this is equal to SSE in linear regression and 1 minus deviance is divided by null devian that is equivalent to multiple R square in linear regression.

So, these are two matrix that could be used to assess the overall fit of the model and then once this is done then we look at the single predictors. We look at whether they are significant or not as I talked about and whether we can strike a balance in terms of prediction performance, classification performance, versus the profiling that is and also statistical modelling context.

(Refer Slide Time: 40:04)

## LOGISTIC REGRESSION

- Outcome variable with m classes ( $m > 2$ )
  - Multinomial logistic regression
    - Separate binary logistic regression model for  $m-1$  classes (one class is treated as reference class)
  - Ordinal logistic regression
    - Large no. of ordinal classes: treat ordinal variable as continuous variable and apply multiple linear regression

IIT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE

16

So, with this we move to our next discussion that is about so till now the exercises that we have been doing they were mainly focused on binary classification, and binary logistic regression model we just had 2 classes class 0, and class 1. Can logistic regression we extended to a scenario where we are dealing with more than 2 classes where we are dealing with you know m classes. So, yes it is possible so we will discuss some of those things.

So, first one is multinomial logistic regression so multinomial logistic regression. So, the categories the classes that we have they are you know so the categorical variable is nominal, so, in that case we can apply multinomial logistic regression. So, what happens in multinomial logistic regression first out of those m classes that we have we have to select one we have to pick one as the reference category, and for the remaining  $m - 1$  classes we create separate binary logistic regression model.

So, for each of the  $m - 1$  classes apart from the reference category class. So, we will have  $n - 1$  classes so for each of the  $n - 1$  classes will create separate binary logistic regression model; that means, for a class 1 we will have the scenario where the observation probability of belonging to class 1 and probability of long not belonging to class 1 so that kind of binary scenario we will have and that for each of the  $n - 1$  classes.

So, we will be dealing with  $m - 1$  binary logistic regression equations and so using that we can compute all those probabilities values with the help of the predictors and then the remaining reference class of that probability for that can always be computed by the  $m - 1$  probabilities values for the  $m - 1$  classes.

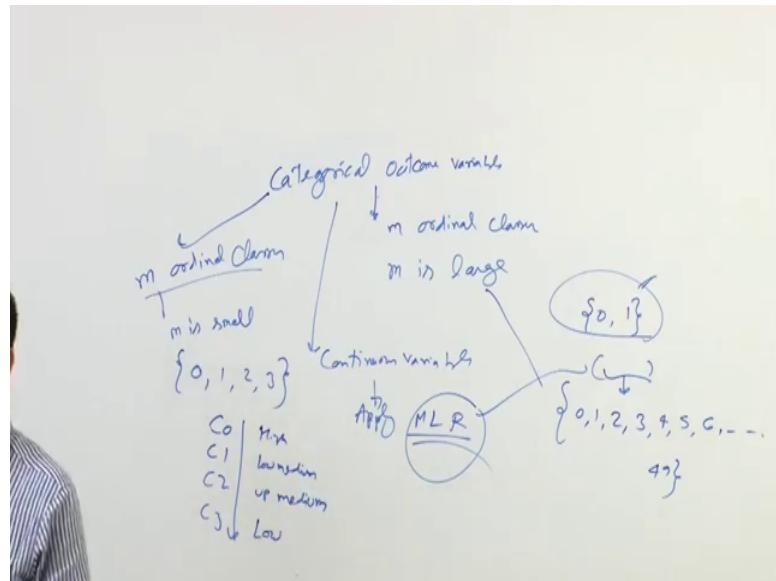
So, we can just subtract that from one for you know and then we can get for the probability value for the different class. And once we have the probabilities values for all  $m$  classes then we can apply our most probable class method routine where the class having highest probability value would be assigned to the new observation.

So, this is how we can go about applying logistic regression to an outcome variable with  $m$  classes. So, this is called this is a multinomial logistic regression and this is applicable mainly to nominal categorical nominal variable right, so the categorical variable having nominal classes.

The second scenario second scenario is about when we have a categorical variable with ordinal classes we have an ordinal variable. So, in those cases we can apply ordinal logistic regression. So, so again within this ordinal logistic regression we can have 2 scenarios. So, as we have understood in some of the initial lectures and supplementary lectures that ordinal variables they have order among different labels is also important that is also meaningful.

So, less than or equal to or greater than or equal to operations they are also applicable in this case, so the first scenario is large number of ordinal classes. So, if our outcome variable which is categorical variable with ordinal classes if that variable have is having large number of ordinal classes, then one solution is treat that ordinal variable as continuous variable and apply multiple linear regression right.

(Refer Slide Time: 44:10)



So, when we have a categorical variable with ordinal classes so we have a categorical variable. Categorical outcome variable when  $m$  ordinal classes, and  $m$  is large then I as I talked about we can treat this particular variable as a continuous variable and apply multiple linear regression.

So, multiple linear regression can be applied and reason is one reason one justification for this is as we talked about earlier that in binary situation we have just two values for a categorical variable and the predicted values up using MLM can range anywhere any real value so, that was the one main problem here.

But when we have  $m$  is large; that means, set of values could be you know many more. So, it could be you know if there are 50 groups let us so in this fashion we can go on up to this. So, the number of values that can be taken by this particular this particular ordinal variable are many more.

So, therefore, the predicted values this can be easily mapped to some of these values could be close to some of these values, and probably multiple linear regression can be still applied so this is one way. When we have a ordinal variable with many number of classes with large number of classes, when  $m$  is large then is still you know instead of logistic regression we can apply multiple linear regression so that is first scenario.

The second one is whatever we have a small number of ordinal classes. So, if m is in this case m ordinal classes that we have if this is a small m is small, then probably we will we will run into the same problem like for binary classification. So, similar problem would be there we will have only few values 0, 1, 2, 3 let us say these many.

So, again small number of classes small number of ordinal classes, so we will have I will run into same problem. So, so we what we do is we use a different version of logistic regression called proportional odds or cumulative logic method as indicated in the slide, so small number of ordinal classes so we would like to use proportional odds, or cumulative cumulative, or logit method.

So, what we do here in this particular method is we create separate binary logistic regression model for m minus 1 cumulative probabilities, so, we talked about that when we so, when we discussed multinomial logistic regression; So, for all m minus 1 classes will have separate binary logistic regression model.

(Refer Slide Time: 47:50)

## LOGISTIC REGRESSION

- Outcome variable with m classes ( $m > 2$ )
  - Ordinal logistic regression
    - Small no. of ordinal classes: Proportional odds or cumulative logit method
      - Separate binary logistic regression model for  $m-1$  cumulative probabilities
- For a three class case: C1, C2, and C3 and a single predictor  $x_1$ 
$$\text{logit}(C1) = \alpha_0 + \beta_1 x_1$$
$$\text{logit}(C1 \text{or } C2) = \beta_0 + \beta_1 x_1$$
- RStudio

17

However, if you see that here will have separate binary logistic regression model not for  $m - 1$  classes not for presence of that class or absence of that class, but  $m - 1$  cumulative probabilities.

So, let us understand what we mean by that so let us take an example for a 3 class case as written in this slide for a 3 class case C 1, C 2, and C 3. Let us say these are the 3 classes

ordinal classes C 1, C 2, C 3 and a single predictor  $x_1$  that is being used. So, our logit equations could be something like this logit for C 1 it could be alpha 0 plus beta 1  $x_1$  and logit for C 1, or C 2 so; that means, from C 1 first logit equation is just for you know observation belonging to C 1. The second is observation belonging to C 1, or C 2 so that gives us the cumulative sense. So, as we talked about that ordinal the order is important so that that means, you know different classes can be compared.

So, therefore, C 1, or C 2 is a you know is a meaningful here in the sense that if we look at the rights part of the equation beta 0 plus beta 1,  $x_1$  you can see that beta 1 is same and both these equation  $x_1$  so you can see because the this is the comparison can be done.

So, the coefficient so the intercept R only difference because the comparison so one is when we talk about ordinal, so one particular class this ordering this ordering is you know this is meaningful. So, when the ordering of classes is meaningful; that means, one can be you know said you know less than one particular class the or higher than one particular class just like the categories that we might have high, you know low medium, upper medium, medium and then low.

So, this kind of this kind of classes we might have all we might have you know a strongly agree to strongly disagree. So, those kind of ordinal classes where the order is meaningful, where the order is meaningful then in those regression we can have something like this cumulative probabilities values, and beta coefficient is can be used the same beta coefficient can be used for both these logistic regression.

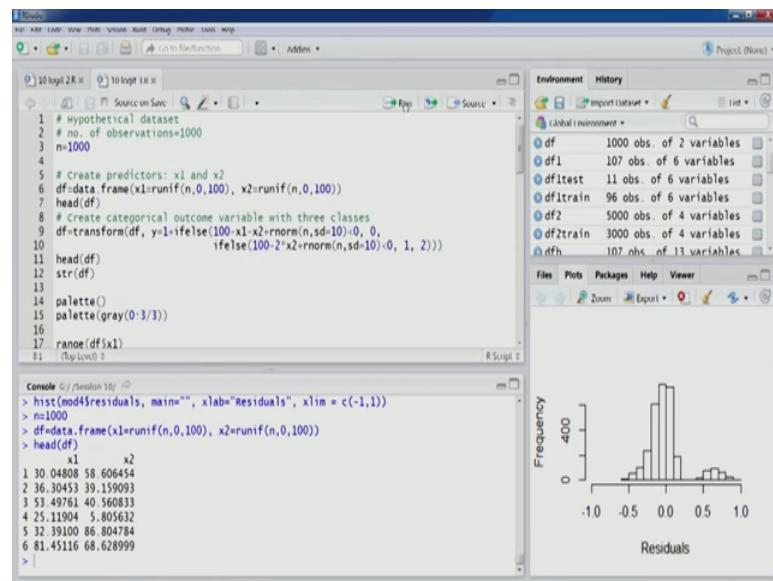
And from this we can compute the cumulative probabilities values and once these cumulative probabilities value values for two of the classes for these classes have been computed the actual probabilities value for C 1, C 2, C 3 can be derived from using the probability value that formulation that we have so, it is from logit, we can derive the probabilities values.

So, once these cumulative so from there once we have the probabilities values for all these classes then we can again apply the most probable class method and also assign the class based on the based on the probability value and again this so in this fashion in this fashion we can apply ordinary logistic regression to a scenario with fewer number of ordinal classes.

So, what we will do we will go through R Studio and do an exercise for this when we have classes more than m is greater than 2, so a logistic regression modelling for classes greater than 2. So, what we will do we will create a hypothetical data set here, so number of observations are 100 in this case, so as you can see and it is 1000.

So, let us create this now what we are going to do is we are going to create a data frame having that data frame where we have 2 variables x 1 and x 2, x 1 as you can see we are using run if for x 1 and x 2 both so the observations. So, the val values would lie between 0 and 100 and n number of values would be created that is 1000 values would be created lying between 0 and 100.

(Refer Slide Time: 52:23)



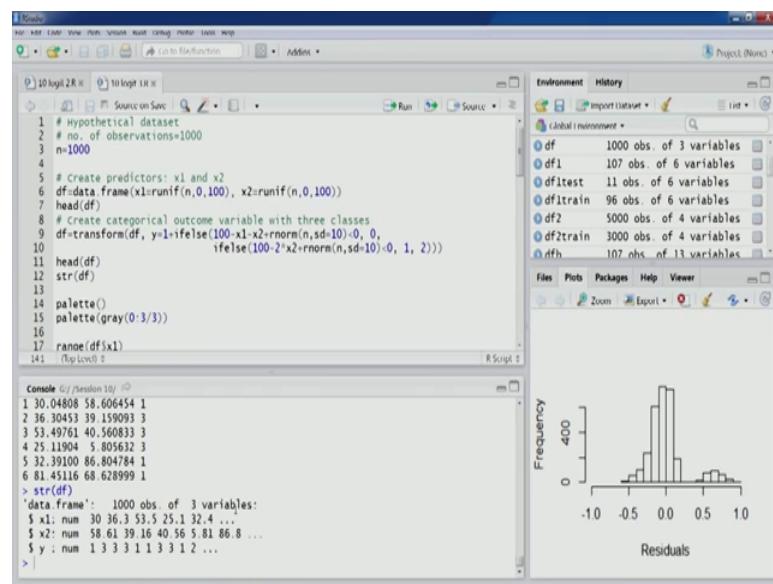
So, let us create this data frame, let us look at some of the observation. For 6 observations you can see all the values they are lying for both the variables x 1 and x 2; they are between 0 and 100.

Now, what we will do we will create a categorical outcome variable with 3 classes. So, this particular data set that we are trying to create we are going to use it for both the scenarios multinomial scenario, and ordinal scenario right, so this is just for illustration purpose so we are not specifying whether the variable is ordinal or not s, we are just going to use it for both the scenarios.

So, what we are going to do is we are using transform function to create our outcome variable, categorical outcome variables. So, you can see that we are trying to compute y as 1 plus if else and if the value is less than 0, this particular value 100 minus x 1 minus x 2 plus again a certain value is being taken from normal distribution no standard deviation 10000 values.

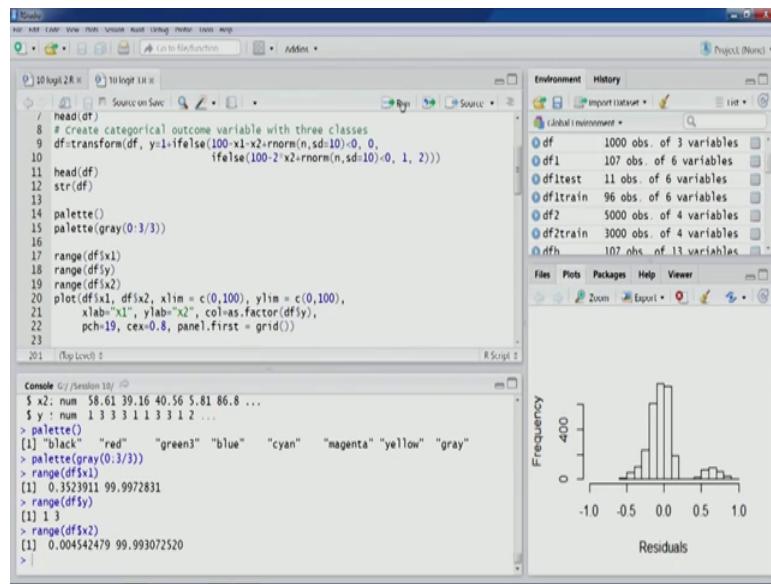
So, if it is less than 0, then you know we are using this information and that is information based on x 1, x 2, and certain additional computation we are trying to assign it a class 0 or then the second if that is not true then second computation so it will get a class 1 or class 2. So, let us compute this once this is done let us look at the observations you can see another variable y that is categorical variable has been created having you can see 1 and 3 2 values are being taken. Let us look at the structure of this particular data frame.

(Refer Slide Time: 54:00)



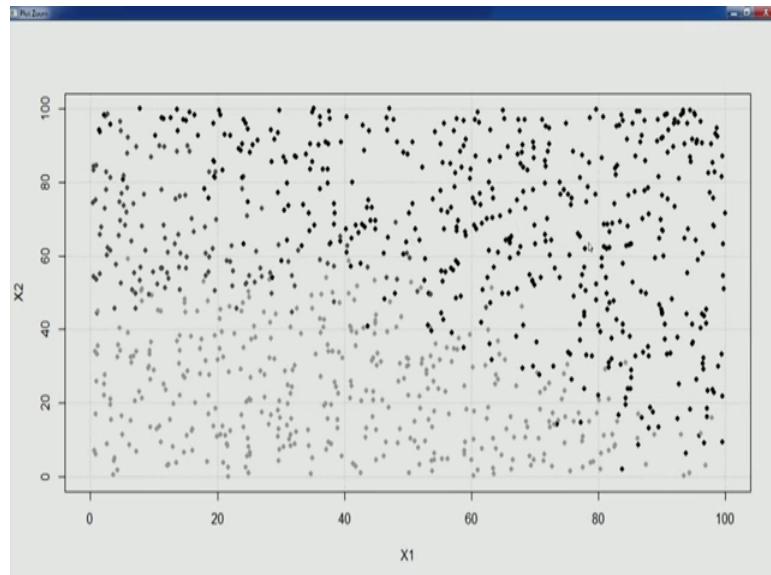
So, these are 3 variables x 1, x 2 values between 0 and 100, and y is taken value 3 values 1, 2, and 3. So, let us plot this particular data set. So, this is our default palette so; however, I would like to use this palette agree 3 sets. So, in this fashion we can have sets of no number of sets that we require. So, let us look at the range of x 1, y, and x 2 which is already.

(Refer Slide Time: 54:25)



Because we have just now created these variables studies actually clearly understood as well so limits 0 to 100, and 0 to 100 and then colouring is with using this particular factor y. So, you can see as dot factor we are using this particular variable and let us create this plot.

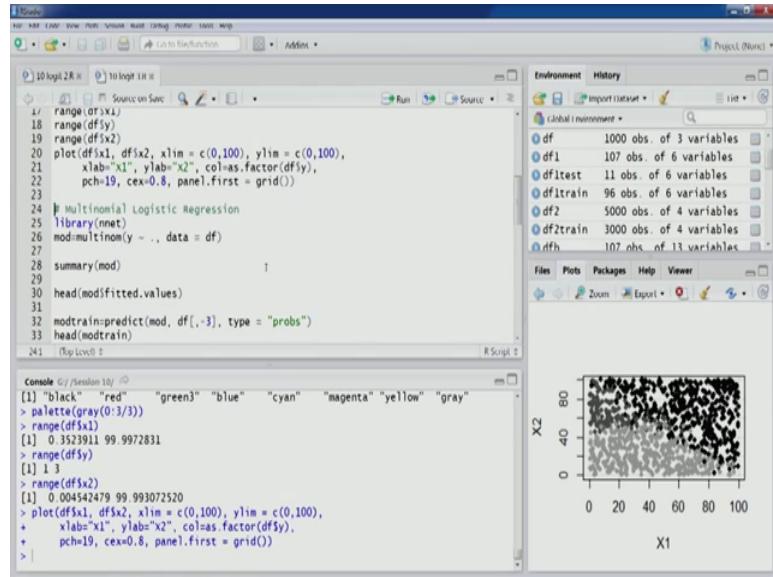
(Refer Slide Time: 54:58)



So, you can see so this is our plotting. So, a one group is here, the second group is having some medium level gray set here and the third group is here, this is lighter gray code

colour. So, these are the values that we are going to use x 1, x 2 this is plot between x 1, and x 2 and the outcome variable has been color coded, so, 3 categories.

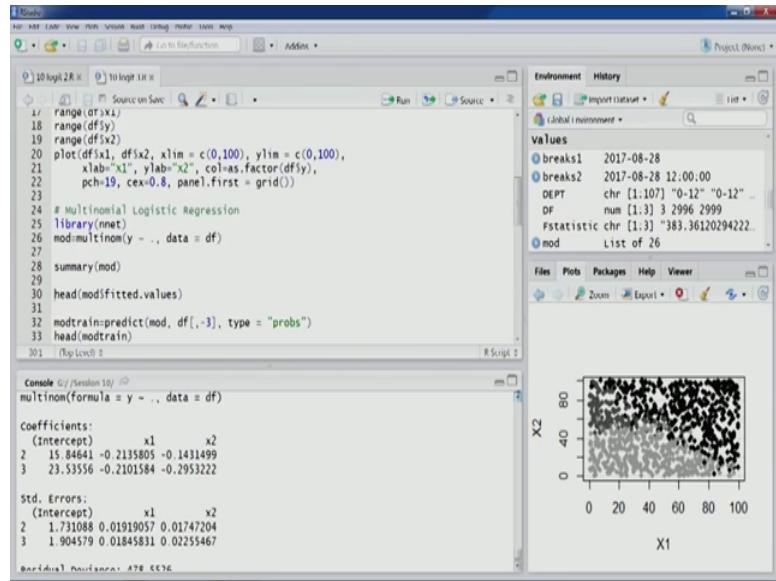
(Refer Slide Time: 55:25)



Now, what we will do is we will use the multinomial logistic regression. So, for this we need this package the library, and net package this package actually for a neural network, but it provides us offers us this function which can be used for multinomial logistic regression.

So, a multi norm is the function so what we are now going to do is the outcome variable y we are going to regress it with the remaining variables that is predictors in this particular data. So, all the observations we are going to use here, so let us run, this let us look at the summary.

(Refer Slide Time: 56:00)



partitions and have probabilities values, estimated probability values, for those new observation.

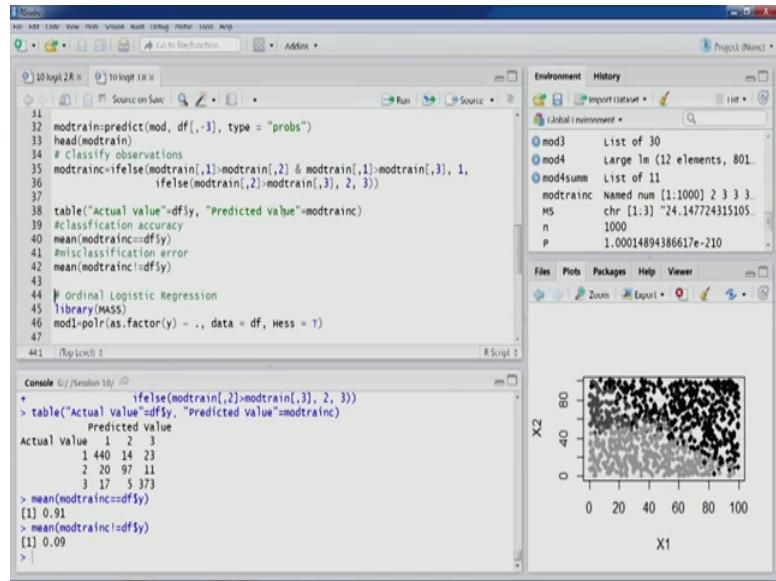
However for the demonstration purpose, we are applying predict function on the training partition that is the full data set itself. So, you can see another argument that we can see is type which is props which is proper probabilities values. So, let us run this and you would see we will have the same probabilities values which were fitted. So, we can see here that same values you can see first row it is same, same values. So, the fitted values and all we have got the same one using predict function.

Now, let us classify these observations. So, this is how we can classify if the probability value for 1 is greater than the probability value for class 2 and also greater than the probability value over class 3, then of course, this observation has to be classified to class 1.

Otherwise we will again compare the probability value for class 2, with probability value for class 3, if again it is greater and I know class 3 value divided then the class 2 is assigned, otherwise class 3.

So, in this fashion we can assign all the observation into appropriate classes so this is implementation of most probable class method. So, once this is done you can see mod train see for all 1000 variables have been created and observations have been assigned to go to appropriate classes.

(Refer Slide Time: 58:13)



Now, let us look at the classification matrix. So, we can see actual value, predicted values, now till now the classification matrices that we have been observing that we have been creating they had 2 values 0 and 1. And we had 2 by 2 classification matrix now this time we are seeing 3 by 3 classification matrix.

So, 3 possibilities for actual values 1, 2, and 3 and the predicted values and we have the corresponding numbers here. So, you can see that again the diagonal element; that means, these 3 elements they represent the values records which have been correctly classified, and the remaining observation off diagonal elements you know they have the records, they have the counts of record, which have been incorrectly classified right.

So, from this we can compute the classification accuracy that is 91 percent in this case, and error that is the remaining 9 percent. So, in this fashion we can apply multinomial logistic regression to a particular data set. Now, let us move to next part that is ordinal logistic regression.

So, in this product case ordinal logistic regression as I talked about 2 scenarios, so 1 where m is large; So, in the in that scenario we can apply multiple regression one exercise that we had done in previous lecture, we had applied multiple linear regression, but that was for the class which had just two you know for a variable category which just had 2 classes 0 and 1.

In the same fashion for the first scenario the multiple linear regression, can be applied. So, there is not much much different in terms of applicability. So, well what we will do we will do an exercise for this scenario where we have to apply this ordinary ordinal logistic regression the cumulative probabilities method, cumulative logit method so for this we need this package mass.

(Refer Slide Time: 60:33)

The screenshot shows the RStudio interface. The left pane displays R code in the console:

```

36 ifelse(modtrain[,2]>modtrain[,3], 2, 3)
37
38 table("Actual value"=df$y, "Predicted value"=modtrain$c)
39 #classification accuracy
40 mean(modtrain$c=df$y)
41 #misclassification error
42 mean(modtrain$c!=df$y)
43
44 # ordinal Logistic Regression
45 library(MASS)
46 mod1=polr(as.factor(y) ~ ., data = df, Hess = T)
47
48 summary(mod1)
49
50 head(mod1$fitted.values)
51
52 modtrain$mod1=mod1$coefficients[,1]
53
46:19

```

The right pane shows the help documentation for the `polr` function from the `MASS` package:

### Ordered Logistic or Probit Regression

**Description**

Fits a logistic or probit regression model to an ordered factor response. The default logistic case is proportional odds logistic regression, after which the function is named.

**Usage**

So, let us load this particular library and we have `polr` is the function for this. So, this is actually for probit however, it can also be used for it can also be used for the cumulative logit method. So, we will just see in the in the help section. So, you can see `polr` this is ordered logistic, or probit regression right.

So, what we are interested in ordered logistic regression, so here ordered logistic regression. So, you can see in the method argument the first one is `draw logistic` this is actually ordered logistic method. This is also called proportional odds logistic regression which we have discussed right.

So, let us go and build this model so `polr` is the function. So, first as you can see `y` variable I have converted into a factor variable and then the request against all other variables that which are predictors data is full, then we need this argument as well as which is mainly if we want to apply for `summary` function later on the model object, so, this is also true.

So, now let us apply summary, so we will get the results we can see the coefficient values here x 1, and x 2, the value and the error and the T values are also there. And we have the residual deviance, and AIC value as well for this particular model and we are interested in finding the fitted values so that is also returned by the model. So, let us look at some of these values.

(Refer Slide Time: 62:00)

The screenshot shows the RStudio interface. The left pane displays an R script with the following code:

```

39 #classification accuracy
40 mean(modtrain==df$y)
41 #misclassification error
42 mean(modtrain!=df$y)
43
44 # Ordinal Logistic Regression
45 library(MASS)
46 mod1=polr(as.factor(y) ~ ., data = df, Hess = T)
47
48 summary(mod1)
49
50 head(mod1$fitted.values)
51
52 mod1train=predict(mod1, df[,-3], type = "probs")
53 head(mod1train)
54 # classify observations
55 mod1train<-ifelse(mod1train[,1]>mod1train[,2]&mod1train[,1]>mod1train[,3], 1,
56 2)
57
58 mod1train

```

The right pane shows the Environment tab with objects like mod1, mod3, mod4, mod4summ, modtrain, M5, and n. Below it, the Help tab is open for the 'Ordered Logistic or Probit Regression' function, showing its description and usage. The Console tab shows the results of the code execution:

```

Residual Deviance: 785.2117
AIC: 793.2117
> head(mod1$fitted.values)
   1        2        3
1 0.1188775524 0.4415837004 0.2395387471
2 0.0515084371 0.2176313067 0.7308584562
3 0.1585359098 0.4442402184 0.2972238718
4 0.0001623156 0.0009184905 0.9988898939
5 0.9697389188 0.0256803680 0.0045807132
6 0.9964164603 0.0030334632 0.0005300764

```

So, you can see for each of these classes class 1, 2, and 3, so these values are actually probability estimated probabilities values. Now using these values we can again apply the most probable class method so first we need to compute the first we need to compute and do the assignment as per the most probable class method like we did in previous exercise.

So, as I talked about predict function can again be used to score new data. So, in this case we are scoring the training partition itself again. So, we expect to get the same values as you can see last row, you can see here, here also the same values are there. So, it is scoring off for the training partition itself. So, we will get the same observation.

Now, what we will do classify observation. So, as we did in previous exercise mod 1 train for class 1 and probability value for class 1 greater than value class 2, and greater than class 3. Then assign it to class 1 otherwise again we look on more comparison the probability value for class 2 is greater than probability value class 3, then assign it to

class 2, otherwise class 3. So, in this way we will have the appropriate classification scores.

(Refer Slide Time: 63:17)

```

48 summary(mod)
49
50 head(mod$fit$values)
51
52 mod$train$predict(mod1, df[,-3], type = "probs")
53 head(mod$train)
54 # classify observations
55 mod$trainc=ifelse(mod$train[,1]>mod$train[,2] & mod$train[,1]>mod$train[,3], 1,
56 ifelse(mod$train[,2]>mod$train[,3], 2, 3))
57
58 table("Actual Value"=df[,y], "Predicted Value"=mod$trainc)
59 #classification accuracy
60 mean(mod$trainc==df[,y])
61 #misclassification error
62 mean(mod$trainc!=df[,y])
63
64 #####
65 | (Top Level) |

```

Console (f:/Session10/)

```

> mod$train$predict(mod1, df[,-3], type = "probs")
> head(mod$train)
  1   2   3
1 0.1188775524 0.4415837004 0.2395387471
2 0.0515084371 0.2176331067 0.7308584562
3 0.2585359098 0.4442402184 0.2972238718
4 0.0001625156 0.0009384905 0.9988699939
5 0.9697389188 0.0256803680 0.0045807132
6 0.9964164603 0.0030534632 0.0005300764
> mod$trainc=ifelse(mod$train[,1]>mod$train[,2] & mod$train[,1]>mod$train[,3], 1,
+ ifelse(mod$train[,2]>mod$train[,3], 2, 3))
> |

```

Now, let us generate the classification matrix here. so you can see 3 by 3 matrix we have 3 classes actual values 3 possibilities, predicted values, 3 predicted classes, so again diagonal elements they represent the correct classification values and off diagonal elements they represent the incorrect classifications.

So, what we can do is so let us compute a classification accuracy, and you can see eighty two percent and the remaining is error.

(Refer Slide Time: 63:49)

The screenshot shows the RStudio interface with the following details:

- File**: R Help, Code, View, Plots, Session, Editor, Photo, Tools, help
- Project**: Project (None)
- Editor**:
  - File: 10 logit.R (active)
  - Code:

```
48 summary(mod1)
49
50 head(mod1fitted.values)
51
52 mod1train<-predict(mod1, df[-3], type = "probs")
53 head(mod1train)
54 # classify observations
55 mod1trainc<-felse(mod1train[,1]>mod1train[,2] & mod1train[,1]>mod1train[,3], 1,
56           ifelse(mod1train[,2]>mod1train[,3], 2, 3))
57
58 table("Actual Value"=df$y, "Predicted Value"=mod1trainc)
59 #classification accuracy
60 mean(mod1trainc==df$y)
61 #misclassification error
62 mean(mod1trainc!=df$y)
63
64 #####
65
```
  - Console: G:/Session 10.R
  - Output:

```
> table("Actual Value"=df$y, "Predicted Value"=mod1trainc)
   Predicted Value
Actual Value      1      2      3
               1 439 31    7
                  2 48 25 55
                  3 10 28 357
> mean(mod1trainc==df$y)
[1] 0.821
> mean(mod1trainc!=df$y)
[1] 0.179
>
```
- Environment**:
  - break2 2017-08-28 12:00:00
  - DEPT chr [1:107] "0-12" "0-12" ...
  - DF num [1:3] 3 2996 2999
  - Statistic chr [1:3] "383.36120294222...
- History**: List of 26
- Files**: Plots, Packages, Help, Viewer
- Help**: Ordered Logit or Probit Regression

So, with this we have completed our discussion on logistic regression and so today we have been also able to cover the scenarios where more than 2 classes are present in our categorical variable. What happens when the classes are nominal and how we can apply logistic regression when classes are ordinal, so we have seen that we have also done an exercise in R.

So, we stop here and we will continue our discussions in next structure for a new technique.

Thank you.