

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

Lecture - 35
Naive Bayes - Part V

Welcome to the Business Analytics and Data Mining Modeling Using R. So, in the previous lecture we were discussing Naive Bayes and specifically we were doing an exercise and modeling exercise in R environment. So, we were able to import the data set and then the process and some of the transformation that we did in the previous lecture. We also looked at some of the issues that we encounter related to dates and arrival and departure time and how we were able to correct those dates that come because of the importing of data from an excel windows, PC to R environment.

And then later on be process to be created the different time intervals based on the departure time, actual departure time. We also did our mod exercise the partitioning and modeling. So, let us, we also looked at the tables the conditional probabilities values for different you know combination outcome variable and predictor combination for different values that is nothing but different categories for the predictors and the outcome variables. So, we looked at all those things.

Now, some of these things can also be performed using a pivot table and excel some of these conditional probabilities. So, though we rely on the computation in R using this particular functions, but these are essentially the mathematical computation computation essentially the technique is more of mathematical in nature and the probability computational or all the crux of this Naive Bayes modeling. So, what we will do? We will export the data set, and that for the training partition that we have done we have created this training partition that we have created in the excel format. So, we will also learn how to export data into excel format.

(Refer Slide Time: 02:12)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Shows an R script named "Rtest.R" with the following code:

```
library(e1071)
mod-naiveBayes(Flight.Status ~ ., dftrain)
attributes(mod)
mod
mod$apriori
mod$tables
# export excel file for pivot table
write.xlsx(dftrain, "D:/roorkee/NPTEL/Session 8/FlightDetails1.xlsx")
mod$tables$Flight.Carrier
mod$tables$Flight.Carrier[1,3]
mod$tables$Flight.Carrier["ontime","Indigo"]

```
- Console:** Shows the output of the R code:

```
Day
Sunday Monday
delayed 0.1666667 0.8333333
ontime 0.1250000 0.8750000
DEPT
DEPT
Y 0-6 6-12 12-18 18-24
delayed 0.2083333 0.6250000 0.1250000 0.0416667
ontime 0.3000000 0.6750000 0.0250000 0.0000000
```
- Environment:** Shows the global environment with objects like `dftest` (43 obs. of 6 variables), `dftrain` (64 obs. of 6 variables), `mod` (List of 4), and `values` (breaks, DEPT, labelsv).
- Help:** Shows the `format` function's usage and arguments.

So, and then we will do a pivot table we will create a pivot table and see how the same conditional probabilities that we have computed using a Naive Bayes function that can also be done using pivot table.

So, let us find and approve it folder for this first. So, if this is the folder then we can specify in this name here and as we have discussed before that we need forward slashes and not the backward slashes in R environment to be able to use the absolute path or files. So, right dot x l s x is the function if we want to export the data into an excel file. So, the training partition data is going to be exported here now we execute this line it has been processed and a file flight details one has been created.

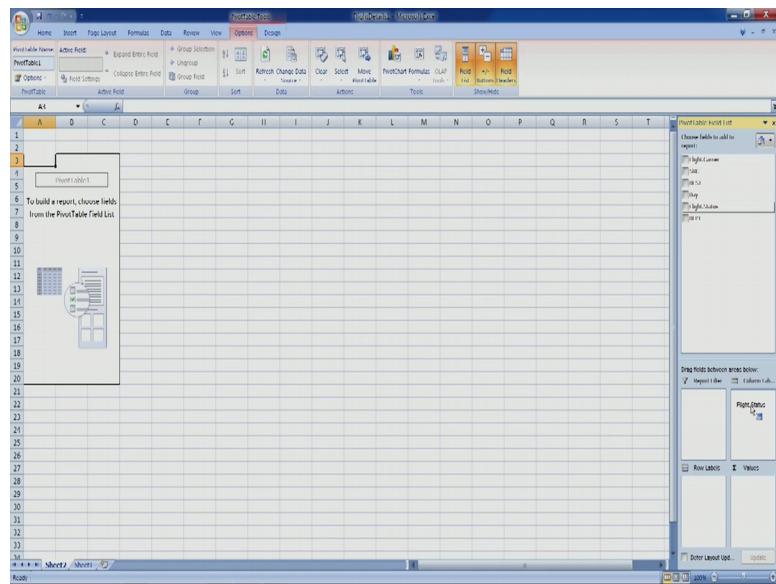
(Refer Slide Time: 03:37)

The screenshot shows a Microsoft Excel spreadsheet titled 'Flight Carrier'. The data is organized into columns A through T, with headers: Flight, Carrier, SRC, DEST, Day, Flight_Status, DEPT. The data consists of 30 rows of flight information. Overlaid on the spreadsheet is the 'Create PivotTable' dialog box. The dialog has three tabs: 'Choose the data that you want to analyze' (selected), 'Choose where you want the pivot table report to be placed', and 'Format Selection'. Under 'Choose the data that you want to analyze', the 'Select a table or range' dropdown shows 'Flight Carrier' (A1:T30). The 'New Worksheet' radio button is selected under 'Choose where you want the pivot table report to be placed'. At the bottom right of the dialog are 'OK' and 'Cancel' buttons.

So, let us open this. So, this is the file you can see first, first column is nothing, but the serial numbers. So, nothing but the index is also these are indices for which were selected when we created the partition. So, right now we do not need this further, column. So, you will delete this and then we have the other columns we have the 6 variables that we require for our modeling exercise. So, these are the 6 variables as you can see here.

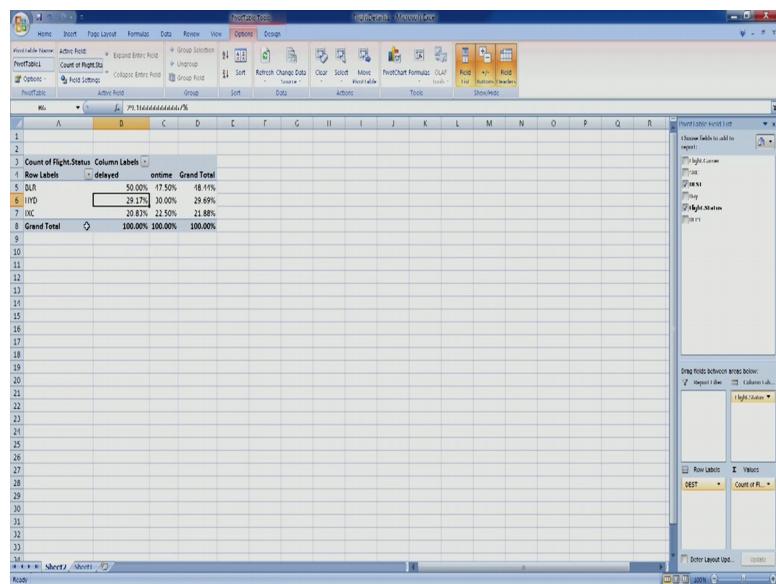
Now, to create the pivot table will select the data set and including the header. So, we will select using ctrl shift and down arrow. So, once this is done we can go to the insert tab within excel and then we can see that pivot table option there and then here we can click on this and we will get the drop down menu and first option is pivot table. So, let us create this. So, the range is already selected.

(Refer Slide Time: 04:53)



Now, we want this pivot table to created on a new worksheet, so we will just do and this is what we get. Now, from here we can create pivot table for different pairs of variables. So, one is flight status this definitely has to go here and column and then for example, depending on the destination or example you want to calculate values for destination and so that we can do.

(Refer Slide Time: 05:14)



So, the column label is flight status row level is destination as you can see a particular pivot table has been created here you can see and the count of count is the default and

has been taken here for flight status few things will change. So, for this left click you would see value field setting and the value field setting. So, everything is ok, and this is summarized by is everything, but show value as there we would like to change and we would like to make it percentage of column. So, this as you can see once we did this. So, the numbers have been converted into percentages.

Now, if we can compare this, to the results that we have got in R to find out whether everything is ok or not, you can see let us look at the results that we had for output a variable and the destination. So, you can see 3 categories b l r, h y d and IXC. So, numbers you can see here 50, 29.17, 20.83 same numbers are here, similarly for on time 47, 47.5, then 30 and then 22.5. So, you can see the probability computation can also be performed using just we just looked at the data created a pivot table and we were able to compute these conditional probabilities. So, the computations that we have performed using the function in R this can actually be performed using a pivot table that are available in excel as well.

So, let us move forward. Now, the table conditional probability table if we are interested in accessing specific values for example, in the table if we just want to have a look at the flight carrier probabilities, conditional probabilities related to flight carrier, we can express in this format because the tables that we have is actually a data frame a list and we can execute this you can see that for a out for outcome variable and the flight carrier we have the conditional probabilities.

(Refer Slide Time: 07:45)

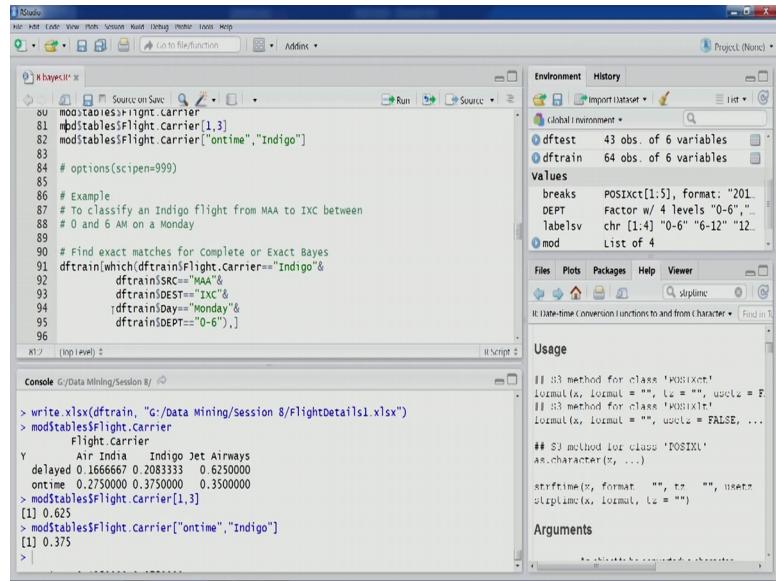
The screenshot shows the RStudio interface with the following details:

- Code Editor (R Script):** Contains R code for creating a Naive Bayes model and exporting it to an Excel file.
- Console:** Shows the execution of the R code, resulting in a 4x4 matrix of values representing probabilities for different flight carriers and on-time status.
- Environment:** Shows variables and their types: `dftrain` (64 obs. of 6 variables), `mod` (List of 4), `DEPT` (Factor w/ 4 levels "0-6", "6-12", "12-18", "18-24"), `labelsv` (chr [1:4] "0-6" "6-12" "12-18" "18-24"), and `breaks` (POSIXct[1:5], format: "201").
- Help:** Shows the `format` function's usage, including its S3 methods for various classes like `POSIXct`, `POSIXlt`, and `POSIXlt`.

If we want just to one first row and third value that is the value for jet airways, so this is how we can you use the brackets notation and axis this value you can see 0.625, if we do not want to use the numbers for row and column we can mention the name of the name of the rows and columns also in this fashion on time and indigo or second example, will get this value 0.375. So, let us execute this you can see. So, this is how we can access the values in R for the Naive Bayes tables from the Naive Bayes tables.

So, let us go through an example where we will try to compute the values by accessing these numbers.

(Refer Slide Time: 08:46)



The screenshot shows the RStudio interface with the following details:

- Code Editor:** Displays R code for a Naive Bayes classifier. It includes filtering for 'Flight Carrier' values like 'Indigo', setting 'DEPT' to '0-6', and specifying 'Day' as 'Monday'. It also includes a 'write.xlsx' command to save the data.
- Console:** Shows the output of the R code, including the resulting data frame 'dftrain' which has 64 observations and 6 variables.
- Environment:** Shows global variables like 'dftest' (43 obs), 'dftrain' (64 obs), and 'mod' (a list of 4).
- Help:** Shows the 'format' function documentation, including its usage, arguments, and examples.

So, probability numbers as we have shown we can also be computed using excel and R as well we can access the individual probabilities value. Now, using them using these value we can now go through one example where in we will try to compute the probability values. So, this is an example.

We want to classify an indigo flight from MAA this particular airport to IXC this particular airport between 0 and 6 am on a Monday. So, you can see in this example. So, all the information related to different predictors is available. Now, using this predictors this information how we can go ahead. As you know that in complete exact way this will have to find the exact matches, which we are not going we are going to use the Naive Bayes computations.

First let us see whether there are any exact matches or not. So, we will just look at whether there are exact matches or not. So, this is how we can do this in the d f train partition and here we can look at whether the flight dot carrier is having this value indigo and that in the source and destination and then day and departure. So, appropriately we can mention the value and find out the rows you can see this particular expression has been written in the row value right, the column value all columns are selected. So, let us execute this, you can see 0 rows. So, no such observation is there in the training partition right. If we just changed this particular value from 0 6 to 6 and 12 let us see if there are any values.

(Refer Slide Time: 10:41)

The screenshot shows the RStudio interface with the following components:

- Script Editor:** Displays the `bayes.R` script containing R code for building a Naive Bayes model.
- Environment Browser:** Shows the global environment with objects like `dftest`, `dfrtrain`, `values`, and `mod`.
- Console:** Shows the output of running the script, including the creation of a `Flight.Status` factor and the resulting data frame `dfrtrain`.
- Plots:** A small plot window is visible in the bottom right corner.

So, we just change the interval you can see there is one value which is actually having all these. So, there are other predictors information or net we have a match for other pattern information except for the departure time interval. So, for that particular interval we did not have any matches therefore, compute exact this cannot be applied.

(Refer Slide Time: 11:14)

The screenshot shows the RStudio interface with the following details:

- Top Bar:** File, Edit, Code, View, Plots, Session, Run, Debug, Profile, Tools, Help.
- Left Panel:** A code editor window titled "R Scripts" containing R code for a Naive Bayes model. The code includes imports for dplyr, magrittr, and rpart, and defines functions for calculating joint probabilities and implementing the Naive Bayes formula.
- Right Panel:** The "Environment" tab is selected, showing global variables like dftrain, dftest, and mod, along with their values and types (e.g., Factor w/ 4 levels).
- Bottom Panel:** The "Console" tab is active, showing the execution of the R code and the resulting output. The output shows the creation of a data frame "dftrain" with columns Flight.CARRIER, SRC, DEST, Day, Flight.Status, and DEPT, and a summary of the data.

So, let us come to the Naive Bayes formula. So, we are going to compute numerator value and denominator value only first. So, as we have discussed in excel exercise as well first we need to compute the probability value for delayed given the information given

the predictor information as per the example. So, you can see is how we can access. So, first a proportion of values belonging to delayed class. So, this is their then delayed and indigo and delayed MAA and delayed IXC delayed 0 to 6 time interval departure time interval and then the later on Monday. So, for all these in this fashion we can actually access the individual condition probability values and we can multiply and the proportion is all is also there. So, this is how we can actually compute the value you can see p 1 has been computed.

(Refer Slide Time: 12:06)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** A script named "K bayes.R" containing R code for calculating the Naive Bayes formula. The code includes filtering for specific conditions like "Monday & dtrain\$DEPT == "0-6"]" and calculating proportions for delayed flights across various categories (Carrier, Destination, Departure Time Interval, and Day).
- Console:** Shows the execution of the script. It starts with the definition of `dtrain` and then proceeds to calculate proportions for delayed flights. The output shows the calculation of `p1` (probability of delayed given the specified conditions) as 0.000706425419560185.
- Environment:** Shows the global environment with variables like `dtrain` (64 observations), `mod` (a list of 4 elements), and `p1` (0.000706425419560185).
- Help:** A search bar and help documentation for R functions.

So, let us print up two values in different digits. So, this is you will get 0.0007064 this is the probability value. Now, for the on time class also we can perform the similar computation here also as you can see that mod and then you can see on time proportion of records are belonging to the on time class. So, that is there.

Then for flight carrier you can see on time and indigo, then for source you can see one time and the MAA this particular airport and then the one on time IXC this the destination airport then the time interval departure time interval this is for on time 0 to 6 and the last value that is for Monday. So, let us compute this. Let us also open the value up to five significant digits. So, now, these are the values for the numerator part.

(Refer Slide Time: 13:09)

```

108 p2<-mod$apriori[["ontime"]]/nrow(dftrain)
109 # (modifiables$flight_carrier["ontime","Indigo"])
110 # (modifiables$SRC["ontime","MAA"])
111 # (modifiables$DEST["ontime","TIC"])
112 # (modifiables$DEPT["ontime","0-6"])
113 # (modifiables$Day["ontime","Monday"])
114 print(p2,digits = 5)
115
116 # Actual probabilities
117 # P(delayed|Example)
118 p1<(p1+p2)
119 # P(ontime|Example)
120 p2/(p1+p2)
121
122 # Scoring test partition
123 modtest$predict(mod, dftest[, -5], type = "class")
124 modtest$predict(mod, dftest[, -5], type = "raw")
1161 [1] 0.0007064

```

Console C:/Data Mining/Session 8/ ↵

```

+ (modifiables$Day["delayed","Monday"])
> print(p1,digits = 4)
[1] 0.0007064
> p2<-mod$apriori[["ontime"]]/nrow(dftrain)
+ (modifiables$flight_carrier["ontime","Indigo"])
+ (modifiables$SRC["ontime","MAA"])
+ (modifiables$DEST["ontime","TIC"])
+ (modifiables$DEPT["ontime","0-6"])
+ (modifiables$Day["ontime","Monday"])
> print(p2,digits = 5)
[1] 0.0034607
>

```

So, as we have done in, as we have discussed before also that we want to compute the actual probability value we can do so, using this p1 we can divide by the summation of these two value we have just computed. So, we will get the actual probability value.

(Refer Slide Time: 13:20)

```

109 # (modifiables$flight_carrier["ontime","Indigo"])
110 # (modifiables$SRC["ontime","MAA"])
111 # (modifiables$DEST["ontime","TIC"])
112 # (modifiables$DEPT["ontime","0-6"])
113 # (modifiables$Day["ontime","Monday"])
114 print(p2,digits = 5)
115
116 # Actual probabilities
117 # P(delayed|Example)
118 p1<(p1+p2)
119 # P(ontime|Example)
120 p2/(p1+p2)
121
122 # Scoring test partition
123 modtest$predict(mod, dftest[, -5], type = "class")
124 modtest$predict(mod, dftest[, -5], type = "raw")
125
1171 [1] 0.0034607
> p1<(p1+p2)
[1] 0.1695237
> p2/(p1+p2)
[1] 0.8304763
>

```

So, these are the actual probabilities value. Now, once we have the probabilities value depending on whether we want to use the most probable class method or depending on whether we want to use the you know we have a class of interest. So, then in that case we

would like to specify the cut off value and compare using that, based on that we can always classify the observation.

(Refer Slide Time: 13:52)

The screenshot shows the RStudio interface with the following components:

- Script Editor:** Displays the R script code.
- Environment View:** Shows the global environment with variables like `modtest`, `p1`, and `p2`.
- Console:** Shows the output of the R code execution.

Script Editor Content:

```
114 print(p2,digits = 5)
115
116 # Actual probabilities
117 # P(delayed|Example)
118 p1/(p1+p2)
119 # P(ontime|Example)
120 p2/(p1+p2)
121
122 # Scoring test partition
123 modest=predict(mod, dftest[,-5], type = "class")
124 modestp=predict(mod, dftest[,-5], type = "raw")
125
126 table("Actual Class=dftest$Flight.status, "Predicted Class"=modtest)
127 head(data.frame("Predicted Class"=modtest,
128 "Actual Class=dftest$Flight.status,
129 "Prob for lSuccess)=modtest[, "delayed"]).
130 dftest[, -5])
131
```

Environment View:

values	mod	p1	p2
breaks	POSIXct[1:5], format: "201		
DEPT	Factor w/ 4 levels "0-6",...		
labelsv	chr [1:4] "0-6" "6-12" "12...		
mod	List of 4		
p1	0.00706425419560185		
p2	0.003460693359375		

Console Output:

```
[1] 0.0034607
> p1/(p1+p2)
[1] 0.1695237
> p2/(p1+p2)
[1] 0.8304763
>
```

Now, once we are done with our building our model using the training partition. So, let us score our test partition and evaluate the performance of this Naive Bayes model on the same. So, what we are going to do here is use the predict function as we have been doing in previous techniques as well and let us say score this.

So, we have two options here type is the argument that can be used within the predict function. So, if we use type as class then we will get the, we will get the scoring for the actual class and if we use this particular time as raw then we will get the estimated probabilities value. So, let us compute each of these two to actual class membership and also the estimated probabilities value.

(Refer Slide Time: 14:41)

The screenshot shows the RStudio interface with several windows open:

- Code Editor:** The main window displays R code for a flight status prediction model. It includes a header section with imports and a body with variable definitions, data frames, and a `modtest` object.
- Environment:** A panel on the right lists global variables and their values, such as `modtest` (a `POSIXct` vector), `breaks` (a character vector), and `DEPT` (a factor).
- Console:** Shows the R command history, including the execution of the code and various R functions like `modtest`, `predict`, and `head`.
- Output:** A panel below the console shows the results of the R commands, such as the head of the data frame and the structure of the `modtest` object.
- Help:** A panel at the bottom provides help documentation for the `strftime` function.

So, let us look at the classification tables that we can generate using table function. So, you can see actual class stored in the flight dot status for test partition as well and the predicted last just be now just now we computed.

(Refer Slide Time: 15:01)

The screenshot shows the RStudio interface with several panes:

- Code Editor:** Displays R code for a Shiny application. The code includes functions for calculating actual probabilities, predicting flight status, and creating a table of actual vs predicted flight status.
- Console:** Shows the output of running the R code, including the calculation of p1/(p1+p2) and the resulting table.
- Environment:** Shows the global environment with variables like modtest, values, breaks, DEPT, labelsv, mod, and modtest.
- History:** Shows a history of previous R commands.
- Plots:** No plots are currently displayed.
- Packages:** Shows available packages: splm, splmDate, and splmDate.
- Help:** Shows help documentation for strftime.
- Viewer:** Shows the resulting table from the R code.
- Usage:** Shows the usage of the strftime function.
- Arguments:** Shows the arguments for the strftime function.

So, these are the results as you can see from this particular classification matrix that 6 delayed flights have been correctly classified as delayed, but 13 of them has been incorrectly classified. On time 4 of them have been incorrectly classified then, but 20 of them have been correctly classified. So, this gives us the idea. So, you can look at 26

observations have been correctly classified and 17 observation have been incorrectly classified this is mainly because of this smaller sample size that is why you are not getting that good result.

(Refer Slide Time: 15:46)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Displays R code for scoring a test partition and calculating classification accuracy.
- Console:** Shows the output of the R code, including a table of results and a summary of the DEPT variable.
- Environment:** Shows the global environment with variables like modtest, dftest, and DEPT.
- Help:** Shows the usage of the strftime function.

```

119 # P(ontime|example)
120 p2/(p1+p2)
121
122 # Scoring test partition
123 modtest<-predict(mod, dftest[, -5], type = "class")
124 modtestp<-predict(mod, dftest[, -5], type = "raw")
125
126 table("Actual Class"=dftest$Flight.Status, "Predicted class"=modtest)
127 head(data.frame("Predicted Class"=modtest,
128                  "Actual Class"=dftest$Flight.Status,
129                  "Prob for 1(success)"=modtestp[, "delayed"],
130                  dftest[, -5]))
131
132 ##Classification accuracy
133 mean(modtest==dftest$Flight.Status)
134 ##Misclassification error
135 mean(modtest!=dftest$Flight.Status)
  
```

	DEPT	delayed	delayed	0.5869109	Jet Airways	BOM	BLR	Monday
8	delayed	delayed	0.6654511	Jet Airways	BOM	BLR	Sunday	
10	delayed	ontime	0.5675214	Jet Airways	BOM	HYD	Monday	
12	ontime	delayed	0.3081360	Air India	BOM	HYD	Monday	
16								

```

DEPT
6 0-6
7 6-12
8 6-12
10 6-12
12 6-12
13 6-12
16 6-12
  
```

Now, we can also have a look at the full information the predicted class the actual class and other on the probability value and the other variables to understand what has happened in the modeling process. As per the results can see predicted class an actual class you can see that these values the probabilities here again.

Again you can see that in this case the default class for our variable default class for our variable we can look at here let us look at the default class you can see delayed and on time. So, the delayed is the going to be the reference category and on time is the main category included in the model. So, would see these probabilities values that p c here is. So, this is a probability of a particular class belonging to on time category, more than 0.5 value of the value is more than 0.5. So, it will belong to the if the does the probability is such then this is going to has been predicted as delayed if it is more than 0.5 if it is less than that and it has been classified as on time.

So, predicted values also we can see here. So, for any probability value that was more than 0.5 we have classified it as delayed and the probability value of less than 0.5 this one from this we can understand it has been classified as on time. So, if we want to look

at the actual classification error accuracy and misclassification number this is all we can do it.

(Refer Slide Time: 18:19)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Contains R code for model evaluation. The code includes:
 - Reading a dataset: `dftest` and `dftrain`.
 - Creating a confusion matrix: `table("Actual Class"=dftest\$Flight.Status, "Predicted class"=modtest)`.
 - Displaying the head of the data frame: `head(data.frame("Predicted Class"=modtest, "Actual Class"=dftest\$Flight.Status, "Prob for 1(success)"=modtest[,5], dftest[,5]))`.
 - Calculating classification accuracy: `mean(modtest==dftest\$Flight.Status)`.
 - Scoring training partition: `modtrain=predict(mod, dftrain[,-5])`.
 - Displaying the head of the training data frame: `table("Actual Class"=dftrain\$Flight.Status, "Predicted Class"=modtrain)`.
 - Calculating classification accuracy: `mean(modtrain==dftrain\$Flight.Status)`.
- Console:** Shows the command history and results. The results for the test partition accuracy are:


```
7 6-12
8 6-12
10 6-12
13 6-12
16 6-12
> str(dftest$Flight.Status)
Factor w/ 2 levels "delayed", "ontime": 2 2 1 2 1 2 1 1 1 1 ...
> mean(modtest==dftest$Flight.Status)
[1] 0.6046512
> mean(modtest!=dftest$Flight.Status)
[1] 0.3953488
>
```
- Environment:** Shows variables defined in the session, including `modtest` (a factor), `mod` (a list), and `dftrain` (a data frame).
- Help:** Shows the help documentation for the `format` function.

So, a comparison between these cold values and the actual values will give us the accuracy. So, that was 0.6 and other was the remaining 0.39 or about 0.4.

(Refer Slide Time: 18:25)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Contains R code for model evaluation. The code includes:
 - Reading a dataset: `dftest` and `dftrain`.
 - Creating a confusion matrix: `table("Actual Class"=dftest\$Flight.Status, "Predicted class"=modtrain)`.
 - Displaying the head of the data frame: `head(data.frame("Predicted Class"=modtrain, "Actual Class"=dftest\$Flight.Status))`.
 - Calculating classification accuracy: `mean(modtrain==dftest\$Flight.Status)`.
 - Scoring training partition: `modtrain=predict(mod, dftrain[,-5])`.
 - Displaying the head of the training data frame: `table("Actual Class"=dftrain\$Flight.Status, "Predicted Class"=modtrain)`.
 - Calculating classification accuracy: `mean(modtrain==dftrain\$Flight.Status)`.
 - Scoring test partition: `cases=1:nrow(dftest)`.
 - Scoring training partition: `modtrain=predict(mod, dftrain[,-5])`.
- Console:** Shows the command history and results. The results for the test partition accuracy are:


```
8 6-12
10 6-12
13 6-12
16 6-12
> str(dftest$Flight.Status)
Factor w/ 2 levels "delayed", "ontime": 2 2 1 2 1 2 1 1 1 1 ...
> mean(modtest==dftest$Flight.Status)
[1] 0.6046512
> mean(modtest!=dftest$Flight.Status)
[1] 0.3953488
> modtrain=predict(mod, dftrain[,-5])
>
```
- Environment:** Shows variables defined in the session, including `modtest` (a factor), `mod` (a list), and `dftrain` (a data frame).
- Help:** Shows the help documentation for the `format` function.

Now, we can also compare the performance of model on test partition with the performance of model on training partition. So, let us score the training partition using the model itself. So, the same function predict.

(Refer Slide Time: 18:43)

The screenshot shows an RStudio interface with the following details:

- Code Editor:** Displays R code for calculating classification accuracy and generating a confusion matrix.
- Console:** Shows the execution of the R code, resulting in a classification matrix table.
- Environment:** Shows variables and their values, including `mod`, `modtest`, and `modtrain`.
- Help:** Shows the `format` function's documentation.

```
130 #classification accuracy
131 mean(modtest==dftest$Flight.Status)
132 #misclassification error
133 mean(modtrain!=dftrain$Flight.Status)
134 #scoring training partition
135 modtrain<-predict(mod, dftrain[,-5])
136 table("Actual Class"=dftrain$Flight.Status, "Predicted Class"=modtrain)
137 #classification accuracy
138 mean(modtrain==dftrain$Flight.Status)
139 #misclassification error
140 mean(modtrain!=dftrain$Flight.Status)
141 # Cumulative Lift Curve
142 cases<-nrow(dftest)
143 cases<-1:nrow(dftest)
144
145 table("Actual Class"=dftest$Flight.Status)
146 > mean(modtrain==dftest$Flight.Status)
[1] 0.6046512
147 > mean(modtest==dftest$Flight.Status)
[1] 0.3953488
148 > modtrain<-predict(mod, dftrain[,-5])
149 > table("Actual Class"=dftrain$Flight.Status, "Predicted Class"=modtrain)
      Predicted Class
Actual Class delayed ontime
  delayed     10    14
  ontime      4    36
> |
```

Predicted Class	Actual Class	Delayed	On Time
Delayed	Delayed	10	14
Delayed	On Time	4	36
On Time	Delayed		
On Time	On Time		

And we can score this we can look at the classification matrix here you can see many more observations are classified correctly in this case as expressed using the accuracy and error numbers.

So, model is performing much better for the training partition main reason being it was built on training partition itself and perform is performing slightly poorly for the test partition. Now, this is also because of the smaller sample size that we have now for this model if we want to generate lift curve.

(Refer Slide Time: 19:19)

The screenshot shows an RStudio interface with the following details:

- Code Editor:** Displays R code for generating a cumulative lift curve.
- Console:** Shows the execution of the R code, resulting in a cumulative lift curve table.
- Environment:** Shows variables and their values, including `mod`, `modtest`, and `modtrain`.
- Help:** Shows the `format` function's documentation.

```
141 mean(modtrain==dftrain$Flight.Status)
142 #misclassification error
143 mean(modtrain!=dftrain$Flight.Status)
144 # Cumulative Lift Curve
145 cases<-nrow(dftest)
146 modtest<-dftest$Flight.Status
147 levels(modtest)<-c("delayed", "ontime")
148 modtest<-as.numeric(as.character(modtest))
149 df1<-data.frame("prob"=modtest[, "delayed"], "actual class"=modtest[, "on time"])
150 df1<-df1[order(-df1$prob),]
151 cumAC<-NULL
152 cumAC[1]<-df1$actual.class[1]
153 for(i in 2:nrow(df1)) {
154   cumAC[i]<-cumAC[i-1]+df1$actual.class[i]
155 }
156
157
158 > mean(modtrain==dftest$Flight.Status)
[1] 0.3953488
159 > modtrain<-predict(mod, dftrain[,-5])
160 > table("Actual Class"=dftest$Flight.Status, "Predicted Class"=modtrain)
      Predicted Class
Actual Class delayed ontime
  delayed     10    14
  ontime      4    36
> mean(modtrain==dftest$Flight.Status)
[1] 0.71875
> mean(modtrain!=dftest$Flight.Status)
[1] 0.28125
> |
```

Predicted Class	Actual Class	Delayed	On Time
Delayed	Delayed	10	14
Delayed	On Time	4	36
On Time	Delayed		
On Time	On Time		

So, this is how we can do it. So, cases we have how many cases. So, for the test partition and I we want to create generate the lift curve cumulative lift curve lift curve. So, we have the 43 observation. So, let us create this.

Now, the values actual classify, actual values for this particular partition test partition let us access them using more test and variable. Now, labels let us change the labels because we need cumulate we need cumulative values for based on the actual, based on the actual classification of the observation. So, we will like to change the labels. So, in this case delayed and on timed one for delayed and 0 for on time.

(Refer Slide Time: 20:18)

The screenshot shows the RStudio interface with the following details:

- Code Editor:** Displays R code for calculating cumulative lift curves. The code includes:
 - Comments: #isclassification error, # Cumulative Lift Curve
 - Variables: modtestn, cases, levels(modtestn), mod, cumAC, df1
 - Operations: mean(modtrain\$Flight.Status), nrow(df1), ifelse(modtestn == 1, 0, 1), cumsum(cumAC), for loops.
 - Print statements: Predicted Class, Actual class delayed ontime, delayed 10 14, ontime 4 36.
 - Final output: mean(modtrain\$Flight.Status) [1] 0.71875, mean(modtrain\$Flight.Status) [1] 0.28125.
- Environment View:** Shows variables and their types:
 - breaks: POSIXct[1:5], format: "%Y-%m-%d %H:%M"
 - cases: int [1:43]; 1 2 3 4 5 6 7 8
 - DEPT: Factor w/ 4 levels "0-6" ...
 - labelsv: chr [1:4] "0-6" "6-12" "12-18" "18-24"
 - mod: List of 4
 - modtest: Factor w/ 2 levels "delayed" "on time"
 - modtestn: Factor w/ 2 levels "1" "0"
- Help View:** Shows the strftime function.
- Console:** Shows the R session history with the same output as the code editor.

So, let us change these labels. Once this is done because this mod test and being the factor variable. So, let us convert it into numeric variable so that we can apply the operations that is the mathematical operations. So, this is you can see now numeric variable and 0 and 1. Now, we would be able to compute the cumulative values right.

So, using this, using the probabilities value that we have that we have computed earlier right for the delayed class right, using this particular value and the actual class that we have just now created a variable we can create a data frame of these two variables.

(Refer Slide Time: 21:05)

```

RStudio
File Edit Window View Plots Session Help Debug Profile Tools Help
Go to file/function Addins
Project (None)
Khadar.R
Source Save Run Source
142 mean(modtrain$Flight.Status)
143 mean(modtrain!=dftrain$Flight.Status)
144
145 # Cumulative Lift Curve
146 cases<-nrow(dftest)
147 modtestn<-dftest$Flight.Status
148 levels(modtestn)<-(1,0) #("delayed", "ontime")
149 modtestn<-as.numeric(as.character(modtestn))
150 df1<-data.frame("prob"=modtestn[, "actual class"=modtestn])
151 df1<-df1[order(-df1$prob),]
152 cumAC=NULL
153 cumAC[1]<-df1$actual.class[1]
154 for(i in 2:nrow(df1)) {
155 cumAC[i]<-cumAC[i-1]+df1$actual.class[i]
156 }
157
15.11 (Top Level) 2
Console C:/Data Mining/Session 8/
> delayed 10 14
> ontime 4 36
> mean(modtrain==dftrain$Flight.Status)
[1] 0.71875
> mean(modtrain!=dftrain$Flight.Status)
[1] 0.28125
> cases<-nrow(dftest)
> modtestn<-dftest$Flight.Status
> levels(modtestn)<-(1,0) #("delayed", "ontime")
> modtestn<-as.numeric(as.character(modtestn))
> df1<-data.frame("prob"=modtestn[, "actual class"=modtestn])
>

```

And then we can sort these variables in terms of the probabilities values the decreasing values right. So, we can if you are interested in looking at the few values. Let us look at the first 6 observation in this particular data frame.

(Refer Slide Time: 21:18)

```

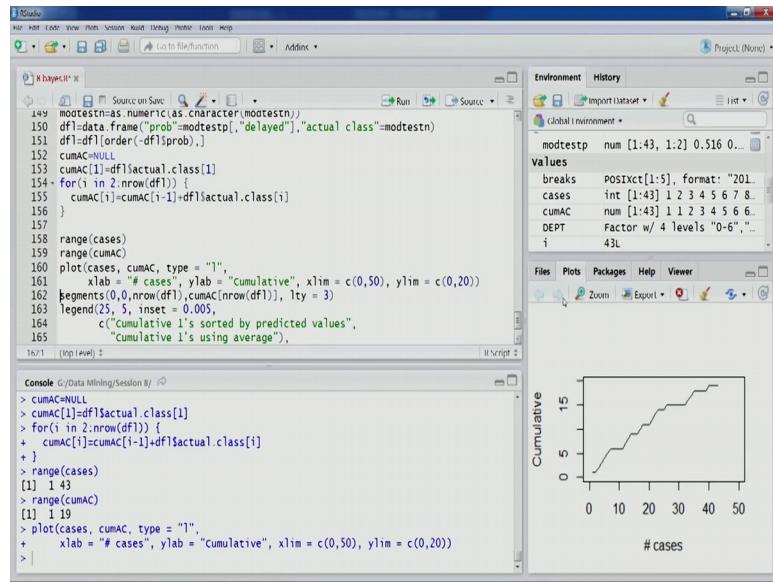
RStudio
File Edit Window View Plots Session Help Debug Profile Tools Help
Go to file/function Addins
Project (None)
Khadar.R
Source Save Run Source
142 mean(modtrain$Flight.Status)
143 mean(modtrain!=dftrain$Flight.Status)
144
145 # Cumulative Lift Curve
146 cases<-nrow(dftest)
147 modtestn<-dftest$Flight.Status
148 levels(modtestn)<-(1,0) #("delayed", "ontime")
149 modtestn<-as.numeric(as.character(modtestn))
150 df1<-data.frame("prob"=modtestn[, "actual class"=modtestn])
151 df1<-df1[order(-df1$prob),]
152 cumAC=NULL
153 cumAC[1]<-df1$actual.class[1]
154 for(i in 2:nrow(df1)) {
155 cumAC[i]<-cumAC[i-1]+df1$actual.class[i]
156 }
157
15.11 (Top Level) 2
Console C:/Data Mining/Session 8/
6 0.3081360 1
> df1<-df1[order(-df1$prob),]
> head(df1)
   prob actual.class
8 0.9524753 1
9 0.9483779 0
10 0.9483779 1
30 0.8245383 1
4 0.6654511 1
2 0.5869109 1
> cumAC=NULL
>

```

You can see probabilities value in actual class, let us order them, as we have done in previous lectures at as well. So, you are again interested in looking at few observations you can see now these values have been sorted or ordered in the decreasing probabilities values sequence right. Now, we can compute the cumulative values. So, let us go through

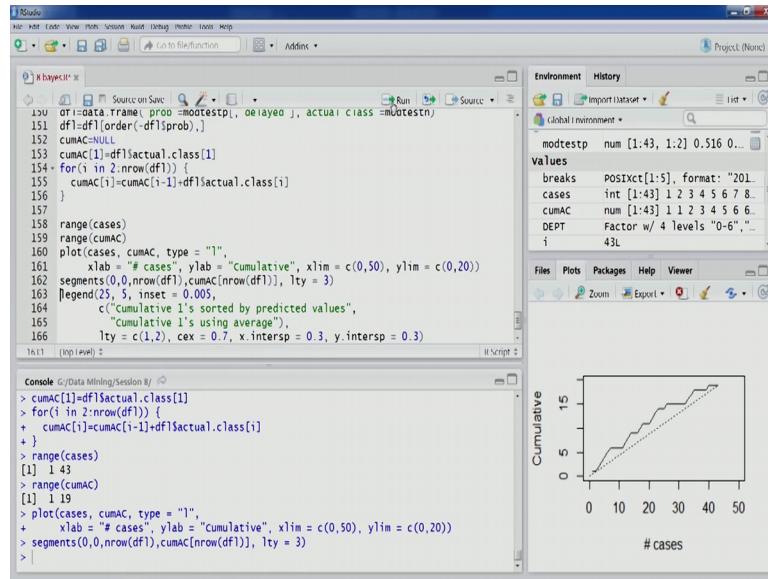
the code for the same. So, the first value is going to be the same. So, how we can do and then for the remaining values we can run this loop when in the previous value is going to be added to get the cumulative value. So, once this is done let us look at the range. So, we have 43 cases, let us look at the cumulative value range 1 to 19 then as you can see appropriately the limits have been specified and other things there.

(Refer Slide Time: 22:19)



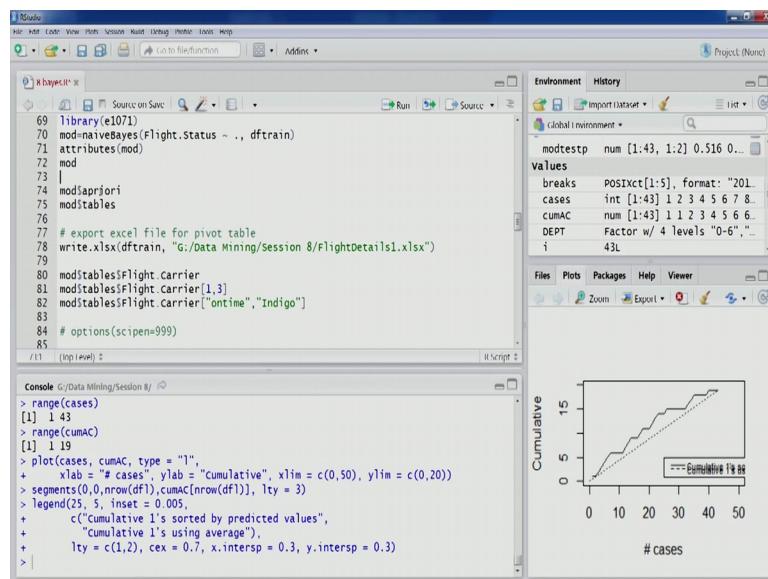
So, we can generate the plot. So, this is the cumulative against curve that we have, this is the plot that we have now if we want to compare its performance with the reference line. So, this is what it is.

(Refer Slide Time: 22:35)



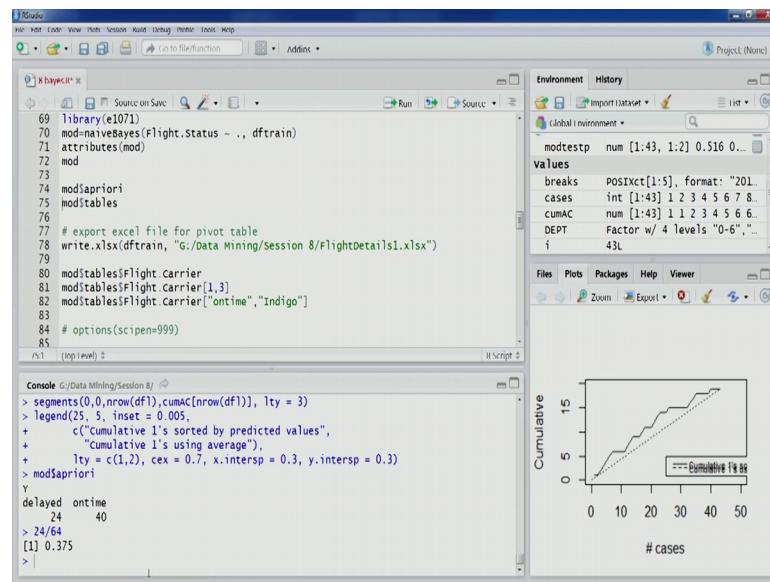
And let us look at the plot. So, you can see this particular plot is performing better than in terms of identifying the most probable ones or most probable delayed flights. So, this model does quite a good job with respect to the average case that is the reference line has so many reference line right. However, if we just look at the prior probabilities and the actual numbers that we have computed right the naive bayes over all on comparing with the probabilities we can find out that if here do the naive rule then what would have been the scenario.

(Refer Slide Time: 23:23)



And if we applied the Naive Bayes modeling, you can see here. So, had we classified all the observations as belonging to the majority class following naive rule then added observation would have been classified as on time and we would have had the error as just 24 divided by 64 and 0.375 right.

(Refer Slide Time: 23:55)



But if we look at the training partition error that we have computed just you know before this right, it is the training partition is model is doing good and training partition we look at the models performance on test partition you see the error is 39.5 and here you is you would seen naive rule the error is 37.5. So, now, rule is actually doing better than our model, but it is the cumulative lift curve that we saw that is giving the usefulness that is indicating towards the usefulness of the model in terms of identifying the most probable ones most probable delayed flights.

So, this is the our modeling exercise and that we wanted perform let us go back to our discussion on Naive Bayes. So, what we will do? We will discuss Naive Bayes few more points.

(Refer Slide Time: 24:57)

NAÏVE BAYES

- Further Comments on Naïve Bayes
 - Good performance despite assumption of independent predictors' values being far from true
 - Requires large no. of records
 - Few classes of predictors might not be represented in the training partition records
 - Zero probability is assumed
 - Good for classification but not for estimating probabilities of class membership

IIT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE

12

So, we look at the Naive Bayes as a technique its quite simple right quite simple and easy to understand is based on the probabilities values. So, the interpretation is quite easy and can give us surprisingly good results in some situations especially when that independent that predictors independent of predictors values that a junction is actually met then it can give us the surprisingly good results or sometimes even outperform other techniques.

So, this good performance can come as mentioned in the first point despite assumption of independent predictors values being far from true. So, even if this assumption is not being met even then we get good performance using Naive Bayes. However, it requires a large number of records because we would like to cover as many scenarios as possible because the predictors they are mainly categorical in nature. So, therefore, the different scenarios that could be there we would like to compute most of them. So, that for a new observation there is always probabilities values for to compute and assign a class to that observation. So, requires a large number of records.

Now, if we do not have a larger sample size then few classes or predictors might not be represented in the training partition record. So, in such situation 0 probability is going to be assumed and then it would be difficult to classify that particular observation right. So however, I mean this particular case is cannot be handled by other techniques as well, but the situation is more complicated in naive Bayes main reason being that the information

in other predictors. For example, for one particular out of 4 or 5 predictors it is just to the one predictor that is not matching, and if that is not if that is then, that is not present in the training partition that one particular class then the value 0 probability value is going to be assumed and therefore, all other predictors. So, the remaining predators information would not be counted which is not the case in other techniques.

Far from this Naive Bayes is mainly suitable for classification tasks as we have discussed in previous lecture previous lectures as well. Another important point is good for a classification, but not for estimating probabilities of class membership. So, as we have talked about that in the Naive Bayes formulation it is the numerator that is sufficient for us to get achieve the rank ordering of rank ordering with, rank ordering of different my classes of outcome variables and that is what we are interested in. So, the model as such will do a good job; however, if we are interested in the actual probabilities values. So, then Naive Bayes is not a suitable technique for such scenarios.

So, this will stop our discussion on Naive Bayes. In the next lecture we will start our discussion on classification and regression trees.

Thank you.