

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 later 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, p c 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)

```
67 dfTest<-df[-partiox,]
68
69 library(e1071)
70 mod=naiveBayes(Flight.Status ~ ., dftrain)
71 attributes(mod)
72 mod
73
74 mod$apriori
75 mod$tables
76
77 # export excel file for pivot table
78 write.xlsx(dftrain, "D:/Roorkee/NPTEL/Session 8/FlightDetails1.xlsx")
79
80 mod$tables$Flight.carrier
81 mod$tables$Flight.carrier[1,3]
82 mod$tables$Flight.carrier["ontime", "Indigo"]
83
```

Console Output:

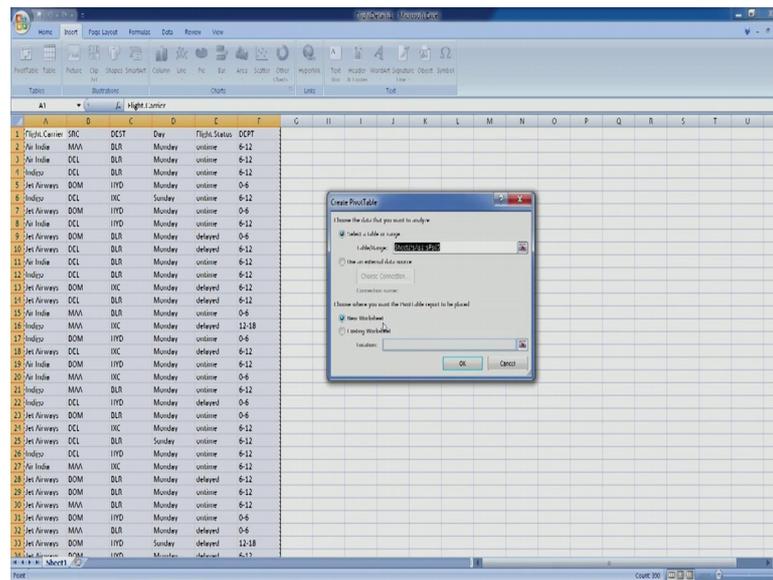
```
Day
Y      Sunday  Monday
delayed 0.1666667 0.8333333
ontime  0.1250000 0.8750000

$DEPT
      DEPT      0-6      6-12      12-18      18-24
Y delayed 0.2083333 0.6250000 0.1250000 0.0416667
ontime  0.3000000 0.6750000 0.0250000 0.0000000
```

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, write 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.

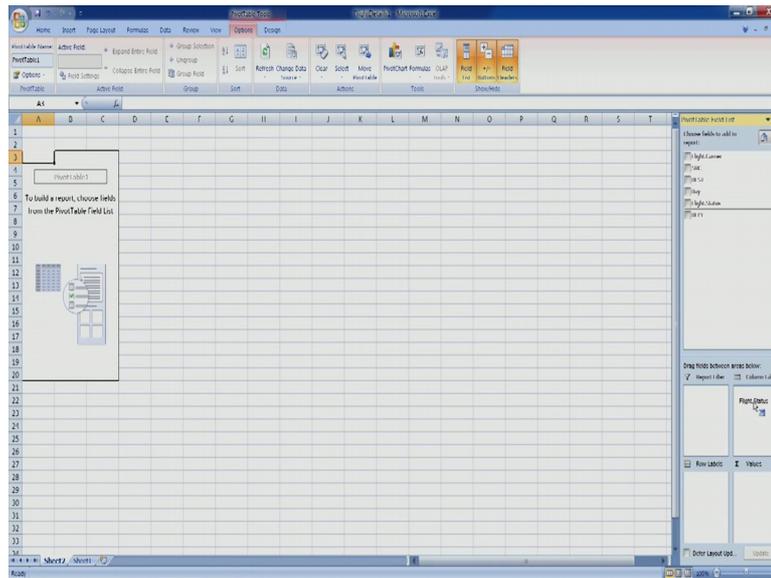
(Refer Slide Time: 03:37)



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.

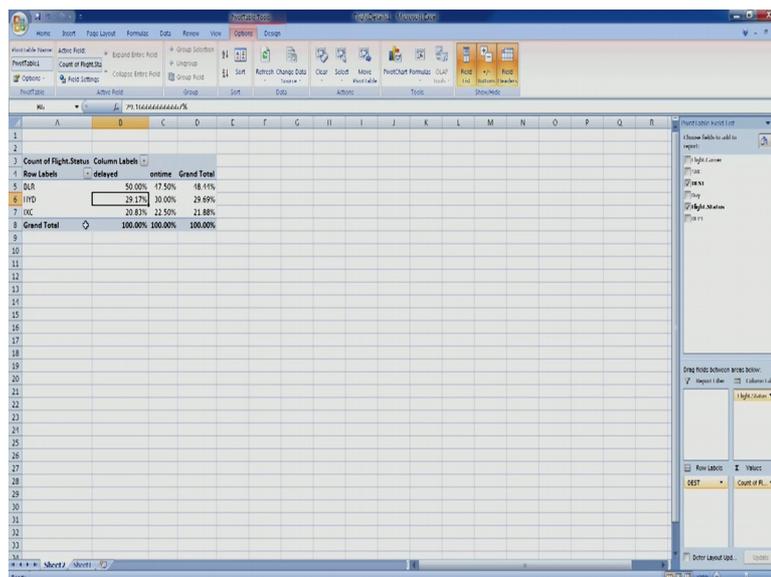
Now, to create the pivot table will select the data set and including the header. So, we will select using ctrl shift and down arrow. So, once this is done we can go to the insert tab within excel and then we can see that pivot table option there and then here we can click on this and we will get the 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 be created on a new worksheet, so we will just do that and this is what we get. Now, from here we can create a 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 brighter 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 brighter 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 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 pub for outcome variable and the flight carrier we have the conditional probabilities.

(Refer Slide Time: 07:45)

```

71 mod$naiveBayes(+flight.status ~ ., dftrain)
72 attributes(mod)
73 mod
74 mod$a priori
75 mod$tables
76
77 # export excel file for pivot table
78 write.xlsx(dftrain, "G:/Data Mining/Session 8/flightDetails1.xlsx")
79
80 mod$tables$flight.carrier
81 mod$tables$flight.carrier[1,3]
82 mod$tables$flight.carrier["ontime", "Indigo"]
83
84 # options(scipen=999)
85
86 # Example

```

```

Y      DEPT      0-6      6-12      12-18      18-24
delayed 0.2083333 0.6250000 0.1250000 0.04166667
ontime  0.3000000 0.6750000 0.0250000 0.00000000

```

```

> write.xlsx(dftrain, "G:/Data Mining/Session 8/flightDetails1.xlsx")
> mod$tables$flight.carrier
      flight.carrier
Y      AIR India  Indigo Jet Airways
delayed 0.1666667 0.2083333  0.6250000
ontime  0.2750000 0.3750000  0.3500000

```

If we want just to one first row and third value that is the value for jet airways, so this is how we again 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 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)

```
80 modStables$flight_carrier
81 modStables$flight_carrier[1,3]
82 modStables$flight_carrier["ontime","Indigo"]
83
84 # options(scipen=999)
85
86 # Example
87 # To classify an indigo flight from MAA to IXC between
88 # 0 and 6 AM on a Monday
89
90 # Find exact matches for Complete or Exact Bayes
91 dftrain[which(dftrain$flight_carrier=="Indigo"&
92             dftrain$src=="MAA"&
93             dftrain$dest=="IXC"&
94             dftrain$day=="Monday"&
95             dftrain$deptime=="0-6"),]
96
```

```
> write.xlsx(dftrain, "G:/Data Mining/Session 8/FlightDetails1.xlsx")
> modStables$flight_carrier
      flight_carrier
Y      Air India  Indigo Jet Airways
delayed 0.1666667 0.2083333 0.6250000
ontime  0.2750000 0.3750000 0.3500000
[1] 0.625
> modStables$flight_carrier["ontime","Indigo"]
[1] 0.375
>
```

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 I 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)

```
84 # options(scipen=999)
85
86 # Example
87 # To classify an Indigo flight from MAA to IXC between
88 # 0 and 6 AM on a Monday
89
90 # Find exact matches for complete or EXACT Bayes
91 dftrain[which(dftrain$flight_carrier=="indigo"&
92             dftrain$src=="MAA"&
93             dftrain$dest=="IXC"&
94             dftrain$day=="Monday"&
95             dftrain$dept=="6-12"),]
96
97 # Naive Bayes formula (Numerator only)
98 # P(delayed|Example)
99 p1=(mod$apriori[["delayed"]]/nrow(dftrain))^
100 (mod$tables$flight_carrier["delayed","indigo"])?
```

```
Console: G:/Data Mining/Session 8/
+ dftrain$dept=="0-6"),]
[1] Flight Carrier SRC DEST Day Flight Status
[6] DEPT
<0 rows> (or 0-length row.names)
> dftrain[which(dftrain$flight_carrier=="indigo"&
+ dftrain$src=="MAA"&
+ dftrain$dest=="IXC"&
+ dftrain$day=="Monday"&
+ dftrain$dept=="6-12"),]
+ Flight Carrier SRC DEST Day Flight Status DEPT
100 Indigo MAA IXC Monday delayed 6-12
>
```

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 except for the departure time interval. So, for that particular interval we did not have any matches therefore, computer exact this cannot be applied.

(Refer Slide Time: 11:14)

```
93 dftrain$dest=="IXC"&
94 dftrain$day=="Monday"&
95 dftrain$dept=="0-6"),]
96
97 # Naive Bayes formula (Numerator only)
98 # P(delayed|Example)
99 p1=(mod$apriori[["delayed"]]/nrow(dftrain))^
100 (mod$tables$flight_carrier["delayed","indigo"])?
101 (mod$tables$src["delayed","MAA"])?
102 (mod$tables$dest["delayed","IXC"])?
103 (mod$tables$dept["delayed","0-6"])?
104 (mod$tables$day["delayed","Monday"])?
105 print(p1,digits = 4)
106
107 # P(ontime|Example)
108 p2=(mod$apriori[["ontime"]]/nrow(dftrain))^
109 (mod$tables$flight_carrier["ontime","indigo"])?
```

```
Console: G:/Data Mining/Session 8/
+ dftrain$dept=="0-6"),]
[1] Flight Carrier SRC DEST Day Flight Status
[6] DEPT
<0 rows> (or 0-length row.names)
> dftrain[which(dftrain$flight_carrier=="indigo"&
+ dftrain$src=="MAA"&
+ dftrain$dest=="IXC"&
+ dftrain$day=="Monday"&
+ dftrain$dept=="6-12"),]
+ Flight Carrier SRC DEST Day Flight Status DEPT
100 Indigo MAA IXC Monday delayed 6-12
>
```

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)

```

94 dftrain$day == Monday &
95 dftrain$DEPT == "0-6",])
96
97 # Naive Bayes formula (Numerator only)
98 # P(delayed|Example)
99 p1=(mod$apriori[["delayed"]]/nrow(dftrain))^
100 (mod$tab$Flight_Carrier["delayed","Indigo"])^
101 (mod$tab$SRC["delayed","MAA"])^
102 (mod$tab$DEST["delayed","IXC"])^
103 (mod$tab$DEPT["delayed","0-6"])^
104 (mod$tab$Day["delayed","Monday"])
105 print(p1,digits = 4)
106
107 # P(ontime|Example)
108 p2=(mod$apriori[["ontime"]]/nrow(dftrain))^
109 (mod$tab$Flight_Carrier["ontime","Indigo"])^
110 (mod$tab$SRC["ontime","MAA"])

```

```

+ dftrain$DEPT=="6-12",])
+ Flight_Carrier SRC DEST Day Flight Status DEPT
+ Indigo MAA IXC Monday delayed 6-12
> p1=(mod$apriori[["delayed"]]/nrow(dftrain))^
+ (mod$tab$Flight_Carrier["delayed","Indigo"])^
+ (mod$tab$SRC["delayed","MAA"])^
+ (mod$tab$DEST["delayed","IXC"])^
+ (mod$tab$DEPT["delayed","0-6"])^
+ (mod$tab$Day["delayed","Monday"])
> print(p1,digits = 4)
[1] 0.0007064

```

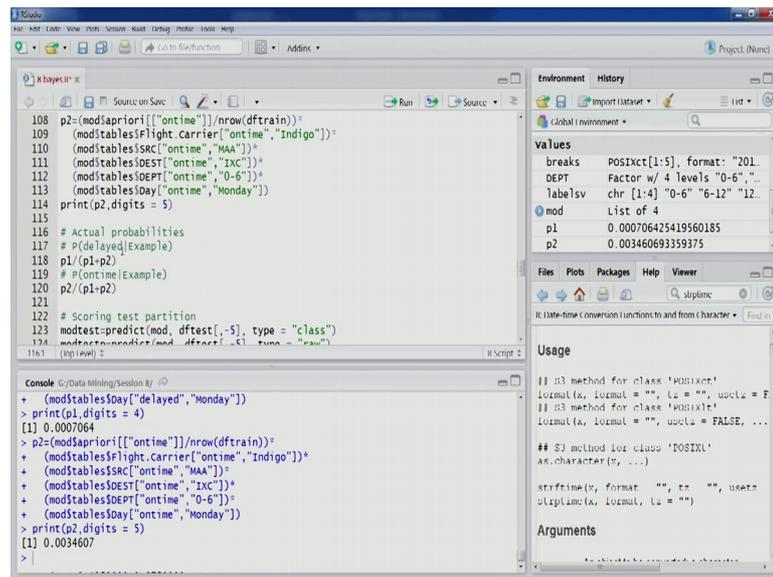
The Environment pane shows the following values:

Variable	Value
dftrain	64 obs. of 6 variables
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	0.000706425419560185

So, let us print up two holes 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 the for the numerator part.

(Refer Slide Time: 13:09)

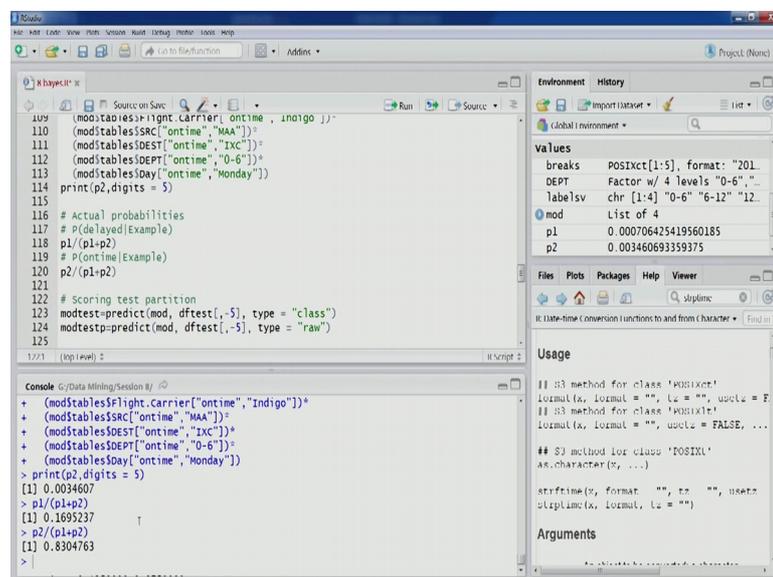


```
108 p2=(mod$apriori[["ontime"]]/nrow(dfttrain))  
109 (mod$tables$flight_carrier["ontime", "Indigo"])*  
110 (mod$tables$src["ontime", "MAA"])*  
111 (mod$tables$dest["ontime", "IXC"])*  
112 (mod$tables$deptime["ontime", "0-6"])*  
113 (mod$tables$day["ontime", "Monday"])  
114 print(p2,digits = 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 modtestp=predict(mod, dftest[,-5], type = "raw")  
125 (top1row) ?
```

```
G:\Data Mining\Session 8/ >  
+ (mod$tables$day["delayed", "Monday"])  
> print(p1, digits = 4)  
[1] 0.0007064  
> p2=(mod$apriori[["ontime"]]/nrow(dfttrain))  
+ (mod$tables$flight_carrier["ontime", "Indigo"])*  
+ (mod$tables$src["ontime", "MAA"])*  
+ (mod$tables$dest["ontime", "IXC"])*  
+ (mod$tables$deptime["ontime", "0-6"])*  
+ (mod$tables$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 p 1 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)



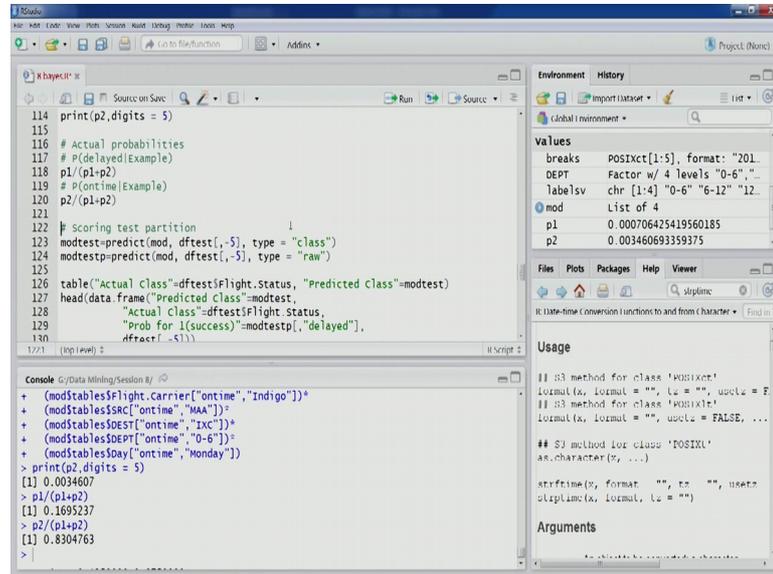
```
109 (mod$tables$flight_carrier["ontime", "Indigo"])*  
110 (mod$tables$src["ontime", "MAA"])*  
111 (mod$tables$dest["ontime", "IXC"])*  
112 (mod$tables$deptime["ontime", "0-6"])*  
113 (mod$tables$day["ontime", "Monday"])  
114 print(p2,digits = 5)  
115  
116 # Actual probabilities  
117 # P(delayed|Example)  
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 modtestp=predict(mod, dftest[,-5], type = "raw")  
125 (top1row) ?
```

```
G:\Data Mining\Session 8/ >  
+ (mod$tables$flight_carrier["ontime", "Indigo"])*  
+ (mod$tables$src["ontime", "MAA"])*  
+ (mod$tables$dest["ontime", "IXC"])*  
+ (mod$tables$deptime["ontime", "0-6"])*  
+ (mod$tables$day["ontime", "Monday"])  
> print(p2, digits = 5)  
[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)



```
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 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               "Prob for 1(failure)"=modtestp[, "ontime"]))
```

Environment History

Global Environment

Values

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	0.000706425419560185
p2	0.003460693359375

Files Plots Packages Help Viewer

Usage

```
## S3 method for class 'POSIXct'
format(x, format = "", tz = "", utc0 = F
## S3 method for class 'POSIXct'
format(x, format = "", utc0 = FALSE, ...
## S3 method for class 'POSIXct'
as.character(x, ...
strftime(x, format = "", tz = "", userz
strftime(x, format = "", tz = ""
```

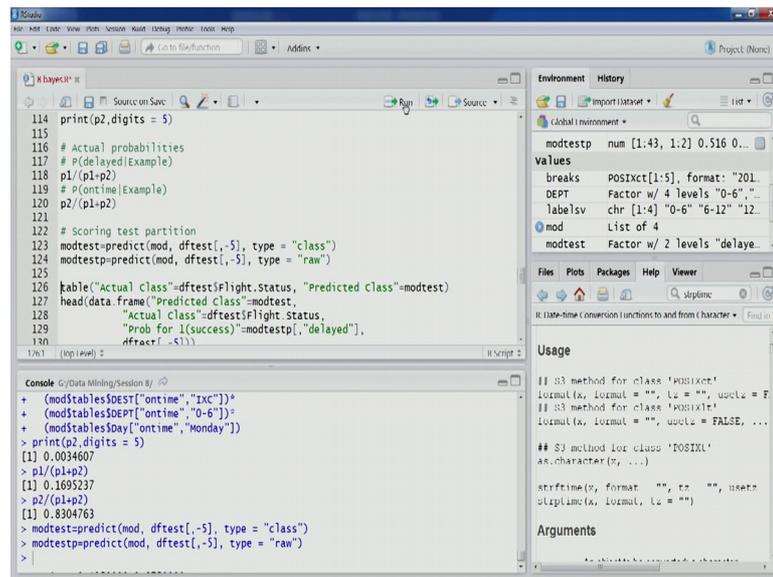
Console

```
+ (mod$tab$Flight.Carrier["ontime", "Indigo"])>
+ (mod$tab$SRC["ontime", "MAA"])>
+ (mod$tab$DEST["ontime", "IXC"])>
+ (mod$tab$DEPT["ontime", "0-6"])>
+ (mod$tab$Day["ontime", "Monday"])
> print(p2, digits = 5)
[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 credit 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)

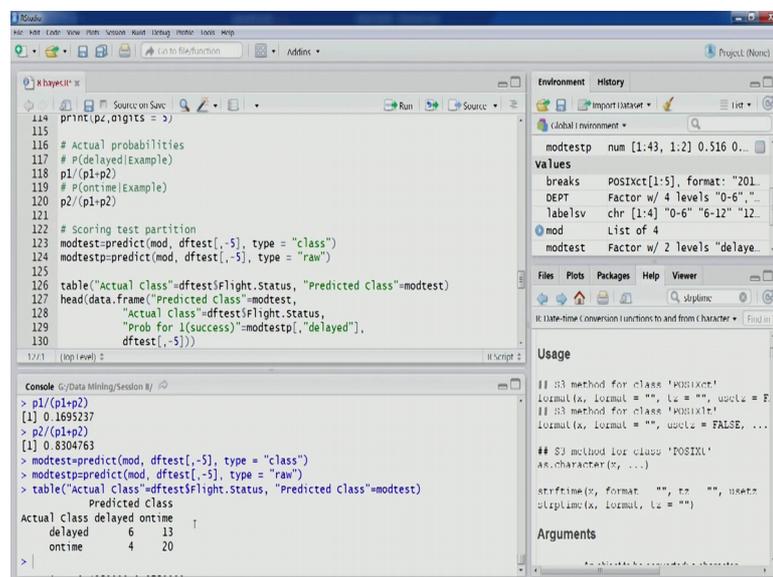


```
114 print(p2,digits = 5)
115
116 # Actual probabilities
117 # P(delayed|Example)
118 p1/(p1+p2)
119 # P(ontime|Example)
120 p2/(p1+p2)
121
122 # Scoring test partition
123 modtest=predict(mod, dftest[,-5], type = "class")
124 modtestpp=predict(mod, dftest[,-5], type = "raw")
125
126 table("Actual 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               dfracr[, -5]))
```

```
> print(p2,digits = 5)
[1] 0.0034607
> p1/(p1+p2)
[1] 0.1695237
> p2/(p1+p2)
[1] 0.8304763
> modtest=predict(mod, dftest[,-5], type = "class")
> modtestpp=predict(mod, dftest[,-5], type = "raw")
>
```

So, let us look at the classification tables that we can generate using table function. So, you can see actual class distort 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)



```
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 modtestpp=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]))
```

```
> p1/(p1+p2)
[1] 0.1695237
> p2/(p1+p2)
[1] 0.8304763
> modtest=predict(mod, dftest[,-5], type = "class")
> modtestpp=predict(mod, dftest[,-5], type = "raw")
> table("Actual Class"=dftest$Flight.Status, "Predicted Class"=modtest)
  Predicted Class
Actual class   delayed   ontime
  delayed      6      13
  ontime       4      20
```

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 observations 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)

```

119 # p1ontime (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
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)"=modtest[, "delayed"],
130               dftest[,-5]))
131
132 #Classification accuracy
133 mean(modtest==dftest$Flight.status)
134 #misclassification error
135 mean(modtest!=dftest$Flight.status)

```

Console Output:

8	delayed	delayed	0.5869109	Jet Airways BOM	BLR Monday
10	delayed	delayed	0.6534511	Jet Airways BOM	BLR Sunday
13	delayed	ontime	0.5675214	Jet Airways BOM	HYD Monday
16	ontime	delayed	0.3081360	Air India BOM	HYD Monday

Environment:

```

modtestp  num [1:43, 1:2] 0.516 0...
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
modtest    Factor w/ 2 levels "delaye...

```

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 properties 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)

```
124 modtestp=predict(mod, dfptest[, -5], type = "raw")
125
126 table("Actual Class"=dfptest$Flight.Status, "Predicted Class"=modtest)
127 head(data.frame("Predicted Class"=modtest,
128 "Actual Class"=dfptest$Flight.Status,
129 "Prob for 1(success)"=modtestp[, "delayed"],
130 dfptest[, -5]))
131
132 #classification accuracy
133 mean(modtest==dfptest$Flight.Status)
134 #misclassification error
135 mean(modtest!=dfptest$Flight.Status)
136
137 # Scoring training partition
138 modtrain=predict(mod, dftrain[, -5])
139 table("Actual Class"=dftrain$Flight.Status, "Predicted Class"=modtrain)
140 #classification accuracy
141
```

```
> str(df$Flight.Status)
Factor w/ 2 levels "delayed", "ontime": 2 2 1 2 1 2 1 1 1 1 ...
> mean(modtest==dfptest$Flight.Status)
[1] 0.6046512
> mean(modtest!=dfptest$Flight.Status)
[1] 0.3953488
>
```

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)

```
131
132 #classification accuracy
133 mean(modtest==dfptest$Flight.Status)
134 #misclassification error
135 mean(modtest!=dfptest$Flight.Status)
136
137 # Scoring training partition
138 modtrain=predict(mod, dftrain[, -5])
139 table("Actual Class"=dftrain$Flight.Status, "Predicted Class"=modtrain)
140 #classification accuracy
141 mean(modtrain=dftrain$Flight.Status)
142 #misclassification error
143 mean(modtrain!=dftrain$Flight.Status)
144
145 # Cumulative Lift Curve
146 cases=1:nrow(dfptest)
```

```
> str(df$Flight.Status)
Factor w/ 2 levels "delayed", "ontime": 2 2 1 2 1 2 1 1 1 1 ...
> mean(modtest==dfptest$Flight.Status)
[1] 0.6046512
> mean(modtest!=dfptest$Flight.Status)
[1] 0.3953488
> modtrain=predict(mod, dftrain[, -5])
>
```

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

(Refer Slide Time: 18:43)

```
130 # Classification accuracy
131
132 #classification accuracy
133 mean(modtest==dftest$Flight.Status)
134 #miscclassification error
135 mean(modtest!=dftest$Flight.Status)
136
137 # Scoring training partition
138 modtrain=predict(mod, dfttrain[,-5])
139 table("Actual Class"=dfttrain$Flight.Status, "Predicted Class"=modtrain)
140 #classification accuracy
141 mean(modtrain==dfttrain$Flight.Status)
142 #miscclassification error
143 mean(modtrain!=dfttrain$Flight.Status)
144
145 # Cumulative Lift Curve
146 cases=1:nrow(dftest)
```

```
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, dfttrain[,-5])
> table("Actual Class"=dfttrain$Flight.Status, "Predicted Class"=modtrain)
      Predicted class
Actual class delayed ontime
delayed      10      14
ontime        4      36
```

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 in my performing much better for the training partition main reason being it was built on training partition itself and perform is performing slightly purely 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 cup.

(Refer Slide Time: 19:19)

```
142 #miscclassification error
143 mean(modtrain!=dfttrain$Flight.Status)
144
145 # Cumulative Lift Curve
146 cases=1:nrow(dftest)
147 modtestn=dftestn$Flight.Status
148 levels(modtestn)=c(1,0) #c("delayed","ontime")
149 modtestn=as.numeric(as.character(modtestn))
150 dfl=data.frame("prob"=modtestp[, "delayed"], "actual class"=modtestn)
151 dfl=dfl[order(-dfl$prob),]
152 cumAC=NULL
153 cumAC[1]=dfl$actual.class[1]
154 for(i in 2:nrow(dfl)){
155   cumAC[i]=cumAC[i-1]+dfl$actual.class[i]
156 }
157
```

```
[1] 0.3953488
> modtrain=predict(mod, dfttrain[,-5])
> table("Actual Class"=dfttrain$Flight.Status, "Predicted Class"=modtrain)
      Predicted class
Actual class delayed ontime
delayed      10      14
ontime        4      36
> mean(modtrain==dfttrain$Flight.Status)
[1] 0.71875
> mean(modtrain!=dfttrain$Flight.Status)
[1] 0.28125
```

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 rest partition let us access them using more test and variable. Now, labels let us change the labels because we need accumulate 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)

```

# Cumulative Lift Curve
cases=1:nrow(dftest)
modtestn=dftest$flight_status
levels(modtestn)=c(1,0) #c("delayed","ontime")
modtestn=as.numeric(as.character(modtestn))
dfl=data.frame("prob"=modtestp[, "delayed"], "actual class"=modtestn)
dfl=dfl[order(-dfl$prob),]
cumAC=NULL
cumAC[1]=dfl$actual.class[1]
for(i in 2:nrow(dfl)) {
  cumAC[i]=cumAC[i-1]+dfl$actual.class[i]
}

```

Console Output:

```

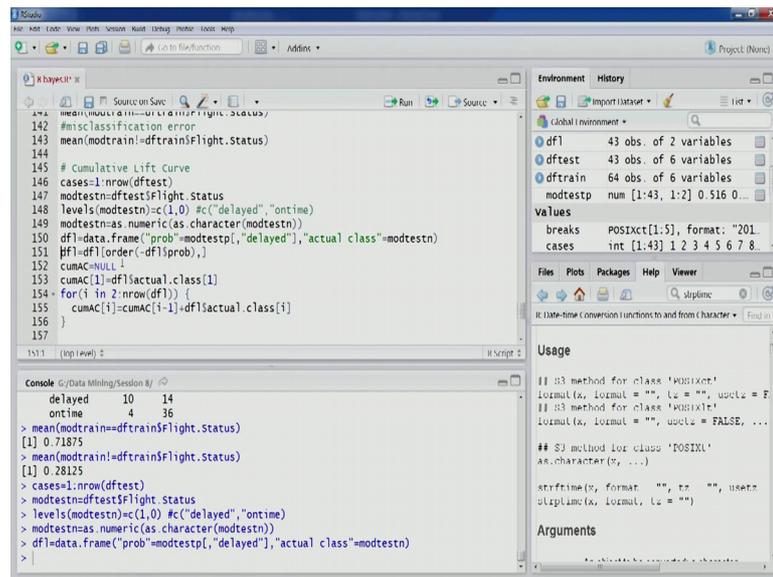
Predicted class
Actual class delayed ontime
delayed 10 14
ontime 4 36
> mean(modtrain=dftrain$flight_status)
[1] 0.71875
> mean(modtrain=dftrain$flight_status)
[1] 0.28125
> cases=1:nrow(dftest)
> modtestn=dftest$flight_status
> levels(modtestn)=c(1,0) #c("delayed","ontime")

```

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)

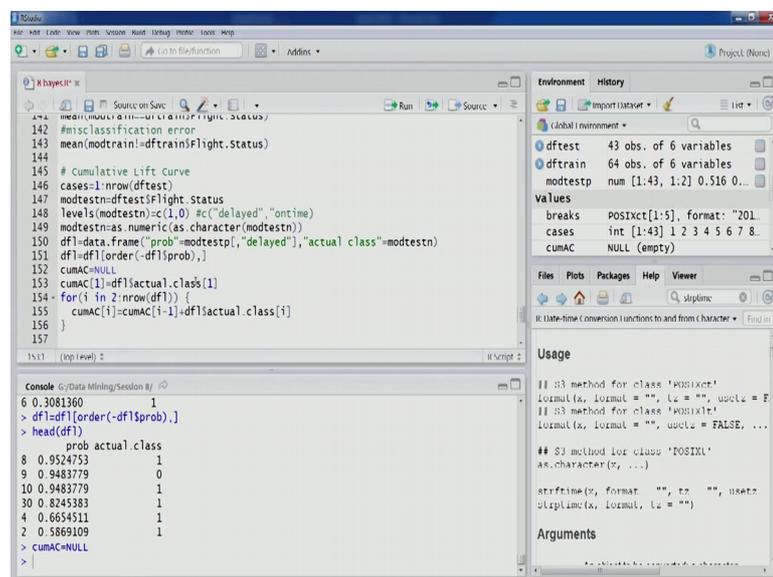


```
142 #miscclassification error
143 mean(modtrain!=dtrain$Flight.Status)
144
145 # Cumulative Lift Curve
146 cases=1:nrow(dftest)
147 modtestn=dftest$Flight.Status
148 levels(modtestn)=c(1,0) #c("delayed","ontime")
149 modtestn=as.numeric(as.character(modtestn))
150 df1=data.frame("prob"=modtestp[,"delayed"],"actual 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
```

```
> mean(modtrain!=dtrain$Flight.Status)
[1] 0.71875
> mean(modtrain!=dtrain$Flight.Status)
[1] 0.28125
> cases=1:nrow(dftest)
> modtestn=dftest$Flight.Status
> levels(modtestn)=c(1,0) #c("delayed","ontime")
> modtestn=as.numeric(as.character(modtestn))
> df1=data.frame("prob"=modtestp[,"delayed"],"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)



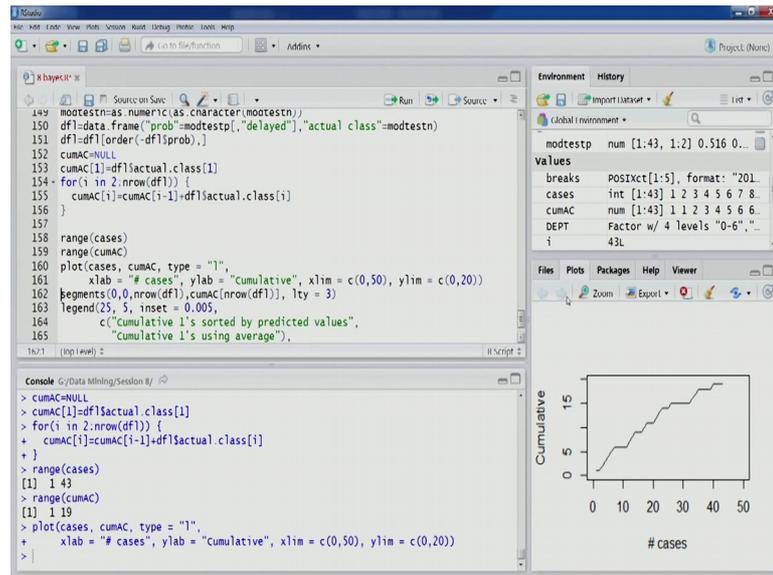
```
150 df1=data.frame("prob"=modtestp[,"delayed"],"actual class"=modtestn)
151 df1=df1[order(-df1$prob),]
152 cumAC=NULL
153 cumAC[1]=df1$actual.class[1]
154 for(i in 2:nrow(df1)) {
155   cumAC[i]=cumAC[i-1]+df1$actual.class[i]
156 }
157
```

```
> 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

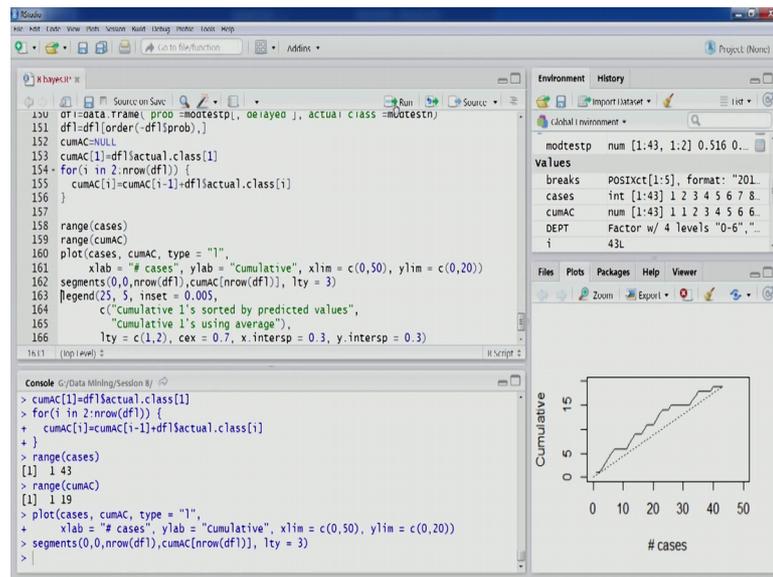
the code for the same. So, the first value is going to be the same. So, how we can do and then for the remaining values we can run this loop when in the previous value is going to be added to get the cumulative value. So, once this is done let us look at the range. So, we have 43 cases, let us look at 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)



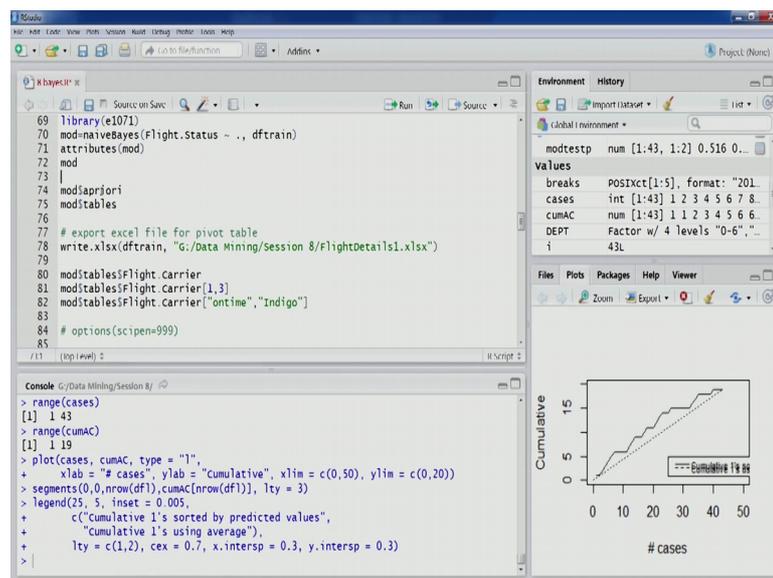
So, we can generate the plot. So, this is the cumulative against curve that we have, this is the plot that we have now if we want to compare its performance with the reference line. So, this is what it is.

(Refer Slide Time: 22:35)



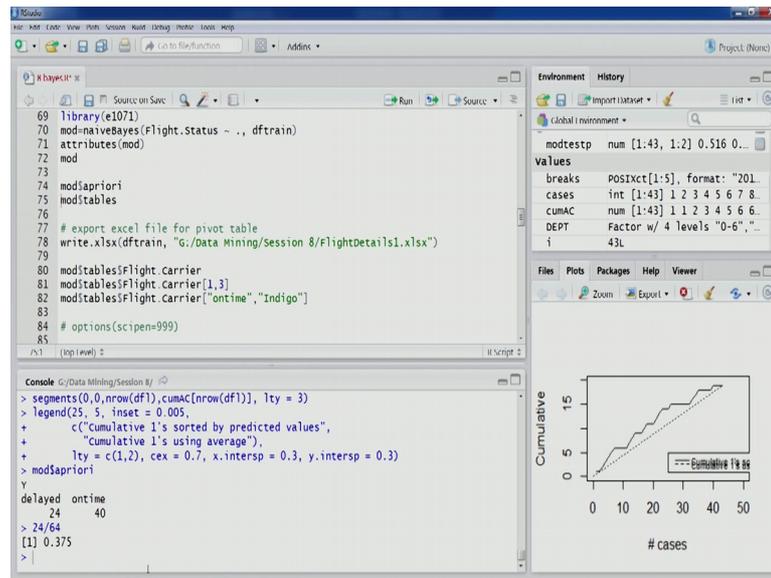
And let us look at the plot. So, you can see this particular plot is performing better than in terms of identifying the most probable ones or most probable delayed flights. So, this model does quite a good job with respect to the average case that is the reference line has so many reference line right. However, if we just look at the prior probabilities and the actual numbers that we have computed right the new base 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 Naive 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 crew. 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 neo base 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 new base 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.