

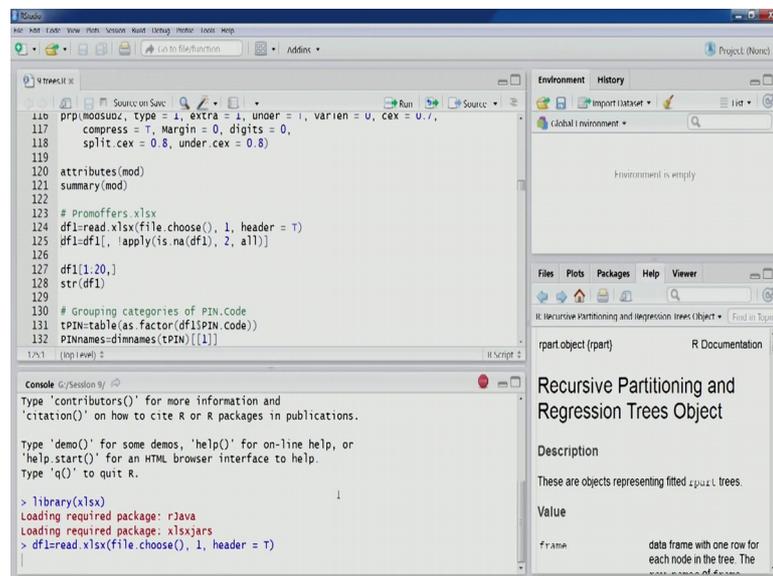
Business Analytics & Data Mining Modeling Using R
Dr. Gaurav Dixit
Department of Management Studies
Indian Institute of Technology, Roorkee

Lecture – 41
Classification and Regression Trees- Part VI

Welcome to the course business analytics and Data Mining Modelling Using R. So, in the previous lecture we were discussing classification trees and specifically we were doing an exercise in R. So, let us start from the point we left in the previous lecture. So, we were we were building classification tree model using the promotional offer data set. So, let us redo few points.

So, that we are able to resume from the same point where we left in the previous lecture. So, let us reload the data set that we used in the previous lecture for promotional offer. So, this is the data set let us import it into the R environment. So, as we saw in the previous lecture this is a quite large data set 5,000 observation. So, it will take slightly more time.

(Refer Slide Time: 01:27)



```
116 ppp=modsuoc, type = 1, extra = 1, under = 1, varien = U, cex = U,,
117 compress = T, Margin = 0, digits = 0,
118 split.cex = 0.8, under.cex = 0.8)
119
120 attributes(mod)
121 summary(mod)
122
123 # Promoffers.xlsx
124 df1=read.xlsx(file.choose(), 1, header = T)
125 df1=df1[, !apply(is.na(df1), 2, all)]
126
127 df1[1:20,]
128 str(df1)
129
130 # Grouping categories of PIN.Code
131 tPIN=table(as.factor(df1$PIN.Code))
132 PINnames=dinames(tPIN)[[1]]
```

Environment History

R Global Environment

rpart object (rpart) R Documentation

Recursive Partitioning and Regression Trees Object

Description

These are objects representing fitted `rpart` trees.

Value

frame data frame with one row for each node in the tree. The

So, you also discussed about the pin code, and how we could? How we were actually using pin code? How we were grouping different categories of pin code? So, let us repeat some of those steps. So, that we are able to reach to the same point. So, pin code promotional offer.

(Refer Slide Time: 02:00)

```
144 y1ad = #Promotional offers Accepted , y1im = c(0,13), cex.names = 0.6)
142 table(as.factor(C_PINcode))
144 # Assign count of PIN Code as its label
145 # PIN Codes having same count will have same label and will be grouped
146 for(x in PINnames) {
147   index=which(as.character(df1$PIN.Code)==x)
148   df1[index,]$PIN.Code=rep(C_PINcode[which(PINnames==x)],length(index))
149 }
150
151 df1$PIN.Code=as.factor(df1$PIN.Code)
152 df1$Promoffer=as.factor(df1$Promoffer)
153 df1$Online=as.factor(df1$Online)
154 tstr(df1)
155
156 # Partitioning: Tr:V:Te>=2500:1500:1000
157 partidx=sample(1:nrow(df1), 2500, replace = F)
```

Environment History

Data

df1 5000 obs. of 9 variables

Values

C_PINcode	int	[1:96]	5	4	5	2	6	3	6	1
index	int	[1:52]	6	56	77	202	238			
PINnames	chr	[1:96]	"110001"	"110000"						
tPIN	'table'	int	[1:96(1d)]	54						

Recursive Partitioning and Regression Trees Object

Description

These are objects representing fitted `rpart` trees.

Value

frame data frame with one row for each node in the tree. The rest is same of `df1`.

Online all these have to be converted fact available as we did in the previous lecture partitioning and then we also build the model.

(Refer Slide Time: 02:11)

```
157 partidx=sample(1:nrow(df1), 2500, replace = F)
158 df1train=df1[partidx,]
159 partidx1=sample(1:nrow(df1)[-partidx], 1500, replace = F)
160 intersect(partidx, partidx1)
161 df1valid=df1[partidx1,]
162 df1test=df1[-c(partidx1, partidx),]
163
164 mod1=rpart(Promoffer ~ ., method = "class", data = df1train,
165           control = rpart.control(cps=0, minsplit = 2, minbucket = 1,
166                                maxcompete = 0, maxsurrogate = 0,
167                                xval = 0),
168           parms = list(split="gini"))
169
170 par(mar=c(0,0,0,0), oma=c(0,0,0,0), xpd=NA)
171 plot(mod1, uniform=T, branch = 0.1, compress = T,
172      margin = 0, nspace = 1)
173 text(mod1, splits = T, use.n = F, all = F, minlength = 0,
```

Environment History

df1test 1000 obs. of 9 variables

df1train 2500 obs. of 9 variables

df1valid 1500 obs. of 9 variables

Values

C_PINcode	int	[1:96]	5	4	5	2	6	3	6	1
index	int	[1:52]	6	56	77	202	238			

mod1 List of 15

Recursive Partitioning and Regression Trees Object

Description

These are objects representing fitted `rpart` trees.

Value

frame data frame with one row for each node in the tree. The rest is same of `df1`.

We need to load this protocol library part library. So, once it is loaded, then we can build the model. We had also saw the tree diagram using a different function the plot function as well as prp function.

(Refer Slide Time: 02:44)

```
1172 margin = 0, nspace = 1)
1173 text(mod1, splits = T, use.n = F, all = F, minlength = 0,
1174       cex=0.7)
1175 prp(mod1, varlen = 0, cex = 0.7, extra = 0, compress = T,
1176      Margin = 0, digits = 0)
1177
1178 # First four levels of full-grown tree
1179 prp(mod1, varlen = 0, cex = 0.7, extra = 0, compress = T,
1180     Margin = 0, digits = 0, nm=T, nn.cex = 0.6)
1181
1182 toss1=as.integer(row.names(mod1$frame))
1183 toss2=sort(toss1)
1184 toss3=toss2[which(toss2==16):length(toss2)]
1185 mod1sub=snip.rpart(mod1, toss = toss3)
1186 prp(mod1sub, varlen = 0, cex = 0.7, extra = 0, compress = T,
1187     Margin = 0, digits = 0)
1188
```

```
Console: G:/Session 9/
+ xval = 0),
+ parms = list(split="gini")
Error in rpart(Promoffer ~ ., method = "class", data = df1train, control = rpart.control(
+ cp = 0, :
+ could not find function "rpart"
+ library(rpart,rpart.control)
+ > mod1=rpart(Promoffer ~ ., method = "class", data = df1train,
+ control = rpart.control(cp=0, minsplit = 2, minbucket = 1,
+ maxcompete = 0, maxsurrogate = 0,
+ xval = 0),
+ parms = list(split="gini"))
+ |
```

Then we started our exercise about finding the split values or finding the split values for different partitions.

So, let us go back. So, this was where we stopped. So, we were calculating the computing the split values.

(Refer Slide Time: 03:07)

```
2214 c1=paste(c1, levels[0:1][as.character(X[j])], sep= ", ")
2215 }
2216 splitvalue[i]=substr(c1, start = 2, stop = nchar(c1))
2217 }
2218 }
2219 } else {
2220 splitvalue[i]=NA
2221 }
2222 }
2223 }
2224 data.frame("NodeNumber"=row.names(mod1$frame), "Splitvar"=mod1$frame$var,
2225           "Splitvalue"=splitvalue, "Cases"=mod1$frame$in,
2226           "Class"=mod1$frame$yval-1,
2227           check.names = F)
2228
```

```
Console: G:/Session 9/
+ }
+ }
+ splitvalue[i]=substr(c1, start = 2, stop = nchar(c1))
+ }
+ }
+ } else {
+ splitvalue[i]=NA
+ }
+ }
+ }
+ |
```

So, this discussion we had already done. So, once for every split variable the variable that were used to develop the full grown tree; once we have computed this split values extracted; rather extracted this split values for all the split variables. So, let us. So, this is

the data frame which we will get give us some useful information. So, the previous the split value that we had computed, now we can see here the same thing which we saw in the.

(Refer Slide Time: 03:42)

```

217 }
218 }
219 else {
220   Splitvalue[i]=NA
221 }
222 i=i+1
223 }
224 data.frame("NodeNumber"=row.names(mod1$frame), "Splitvar"=mod1$frame$var,
225           "Splitvalue"=splitvalue, "Cases"=mod1$frame$n,
226           "Class"=mod1$frame$yval-1,
227           check.names = F)
228
229 mod1train=predict(mod1, dfttrain[,c(3)], type = "class")
230 table("Actual value"=dfttrain$Promoffer, "Predicted value"=mod1train)
231 #classification accuracy
232 mean(mod1train==dfttrain$Promoffer)

```

NodeNumber	Splitvar	Splitvalue	Cases	Class	
1	Income	101.5	2500	0	
2	spending	2.86649178649802	1926	0	
3	4	2.59733373904682	1742	0	
4	8	2.16136573592129	1649	0	
5	<leaf>	<NA>	1440	0	
6	17	2.16641149594615	209	0	
7	34	Income	91.5	208	0

The same thing which we saw in the full grown tree and where we could actually follow the root node and what was the split variable and value combination then the other splits and the; predictor and value combination all right. So, the same thing is represented in this a tabular format in the here and, that is why be required to extract this split value. So, that we could create this tabular format you can see node number. The unique node number as we discussed about the tree diagrams, and they split variable that has been used in that particular node number.

And the corresponding is split value that has been used to create a split and number of cases and the class for that particular node right. So, so, in this fashion we can look at each of the each you know different nodes that are their node number two the split variable was spending. So, this is with respect to full grown tree as we saw in the previous lecture. So, this is the split value for spending and the number of cases that are there in that particular node. So, we want to have a look at the folder diagram again. So, that was quite a.

(Refer Slide Time: 04:56)

```
172 margin = 0, nspace = 1
173 text(modl, splits = T, use.n = F, all = F, minlength = 0,
174     cex=0.7)
175 prp(modl, varlen = 0, cex = 0.7, extra = 0, compress = T,
176     Margin = 0, digits = 0)
177
178 # First four levels of full-grown tree
179 prp(modl, varlen = 0, cex = 0.7, extra = 0, compress = T,
180     Margin = 0, digits = 0, nn=T, nn.cex = 0.6)
181
182 toss1=as.integer(row.names(modl$frame))
183 toss2=sort(toss1)
184 toss3=toss2[which(toss2==16):length(toss2)]
185 modlsub=snip.rpart(modl, toss = toss3)
186 prp(modlsub, varlen = 0, cex = 0.7, extra = 0, compress = T,
187     Margin = 0, digits = 0)
188
```

NodeNumber	splitvar	splitvalue	Cases	Class	
1	Income	101.5	2500	0	
2	spending	2.86649178649802	1926	0	
3	4	2.59733373904682	1742	0	
4	8	2.16136573592129	1649	0	
5	16	<leaf>	<NA>	1440	0
6	17	spending	2.16641149594615	209	0
7	34	Income	91.5	208	0
8	68	<leaf>	<NA>	198	0
9	69	Income	92.5	10	0
10	138	<leaf>	<NA>	9	0
11	139	<leaf>	<NA>	1	1

That was quite a big tree model and so, this will this prp function will give us the tree diagram.

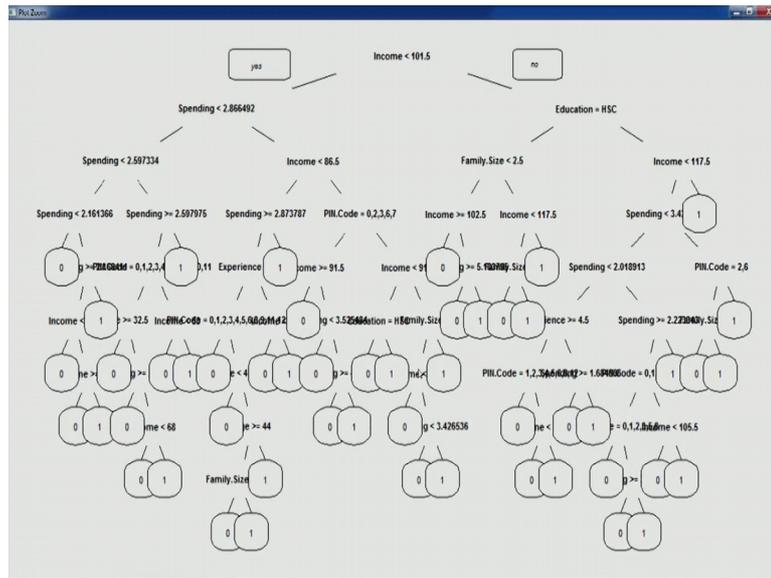
(Refer Slide Time: 05:12)

```
163 control = rpart.control(cp=u, minsplit = 2, minbucket = 1,
164                       maxcompete = 0, maxsurrogate = 0,
165                       xval = 0),
166
167     parms = list(split="gini"))
168
169 par(mar=c(0,0,0,0), oma=c(0,0,0,0), xpd=NA)
170 plot(modl, uniform=T, branch = 0.1, compress = T,
171     margin = 0, nspace = 1)
172 text(modl, splits = T, use.n = F, all = F, minlength = 0,
173     cex=0.7)
174 prp(modl, varlen = 0, cex = 0.7, extra = 0, compress = T,
175     Margin = 0, digits = 0)
176
177 # First four levels of full-grown tree
178 prp(modl, varlen = 0, cex = 0.7, extra = 0, compress = T,
179     Margin = 0, digits = 0, nn=T, nn.cex = 0.6)
180
181
```

```
> prp(modl, varlen = 0, cex = 0.7, extra = 0, compress = T,
+     Margin = 0, digits = 0)
Error in prp(modl, varlen = 0, cex = 0.7, extra = 0, compress = T, Margin = 0, :
could not find function "prp"
> library(rpart.control)
Error in library(rpart.control) :
there is no package called 'rpart.control'
> install.packages("rpart.control")
Warning in install.packages :
cannot open URL 'http://www.stats.ox.ac.uk/pub/Rwin/src/contrib/PACKAGES.rds': HTTP s
tatus was '404 Not Found'
```

So, we need to load this particular library to be able to use the, this particular function; so once this particular package installed will be able to use rpart this function. So, let us reload the library the name was different. So, that was the problem. So, it was actually rpart dot plot. So, prp function is there. So, we have to use this. So, now, this function would be available to us. So, let us plot this you can see this was quite a big plot.

(Refer Slide Time: 06:33)



Tree diagram and the table that we had generated.

(Refer Slide Time: 06:38)

```

163 control = rpart.control(cp=0, minsplit = 4, minbucket = 1,
166 maxcompete = 0, maxsurrogate = 0,
167 xval = 0),
168 parms = list(split="gini"),
169
170 par(mar=c(0,0,0,0), oma=c(0,0,0,0), xpd=NA)
171 plot(mod1, uniform=T, branch = 0.1, compress = T,
172 margin = 0, nspace = 1)
173 text(mod1, splits = T, use.n = F, all = F, minlength = 0,
174 cex=0.7)
175 prp(mod1, varlen = 0, cex = 0.7, extra = 0, compress = T,
176 Margin = 0, digits = 0)
177
178 # First four levels of full-grown tree
179 prp(mod1, varlen = 0, cex = 0.7, extra = 0, compress = T,
180 Margin = 0, digits = 0, nm=T, nn.cex = 0.6)
181

```

NodeNumber	splitvar	SplitValue	Cases	Class
1	Income	101.5	2500	0
2	Spending	2.86649178649202	1926	0
3	Spending	2.59733373904682	1742	0
4	<leaf>	<NA>	1649	0
5	<leaf>	<NA>	1440	0
6	Spending	2.16641149594615	209	0
7	Income	91.5	208	0
8	<leaf>	<NA>	198	0
9	Income	92.5	10	0
10	<leaf>	<NA>	9	0

For split variable and value combination, this is the table that we are generated. Now from here and looking at the tree model you can get a better sense of, what happened in our tree diagram? So, you can see first node, root node income 1, 0 less than 101.5 and, that is the node number and the same value is here you can see in the table as well. Similarly, they say this is the second node expanding. So, the same thing is presented here you can see the value is also same.

So, the tree diagram the information that we have is specifically focusing on the split variable and split value, that we can also present in this tabular format and understand; what are the important predicted value combinations? This is especially important if we have quite a big tree we are developing a full grown tree and. So, in those situation this tabular format might be more useful. And you can see the next node number is 4. So, the 4 is will come here. So, this one is the split variable value combination is spending less than 2.16 right 2.59. So, there must be one mole this is the 1, 4 and then this is the node number 8. So, in this fashion you can see.

So, the; these are this is the table for split variable value combination. So, let us move forward. So, once we have built this model. So, let us look at the performance of this model on training parties itself; since this model has been built you know this has be this 3 model is a full grown tree. So, therefore, what we expect is that? All the observation would be correctly classified, because once we develop full grown tree.

So, the; we keep on creating partitions, we keep on doing splits till we are able to reach till we are able to create pure homogeneous groups or rectangles or partitions right. So, predictor is the function as we have been using in previous techniques as well. So, this can be used to score the a particular data set. So, the model mod 1 and we are again trying to score off the training partition itself and minus c 3.

So, we are not including the; dependent variable outcome variable in this particular data set and the type is class. So, with this we can look at the performance and you would see that, this is the classification matrix you can see actual value and predicted value. So, all the 2251 actual value have been correctly classified as; correctly classified as class 0 and you can see 249 actual class one records they have been correctly classified right.

So, the if we look at the we find the classification accuracy and error you would see it is 100 percent right 1 and 0 here. So, you can see the performance is 100 percent because the tree model is all fitting the data completely. So, the whole data the data set the training partition has been completely fit with 100 percent accuracy and 0 percent error.

Now, let us let us look at the performance of this full grown tree model on validation partition and test partition. So, we had already we have already created these two partitions starts apply the same function.

(Refer Slide Time: 10:15)

```

444 cneck.names = F)
228
229 modtrain=predict(mod1, dfttrain[,c(3)], type = "class")
230 table("Actual value"=dfttrain$Promoffer, "Predicted value"=modtrain)
231 #classification accuracy
232 mean(modtrain==dfttrain$Promoffer)
233 #misclassification error
234 mean(modtrain!=dfttrain$Promoffer)
235
236 modvalid=predict(mod1, dfvalid[,c(3)], type = "class")
237 table("Actual value"=dfvalid$Promoffer, "Predicted value"=modvalid)
238 #classification accuracy
239 mean(modvalid==dfvalid$Promoffer)
240 #misclassification error
241 mean(modvalid!=dfvalid$Promoffer)
242
243 modltest=predict(mod1, dfitest[,c(3)], type = "class")

```

```

labs do not fit even at cex 0.15, there may be some overplotting
> modtrain=predict(mod1, dfttrain[,c(3)], type = "class")
> table("Actual value"=dfttrain$Promoffer, "Predicted value"=modtrain)
      Predicted value
Actual value 0 1
             0 2251 0
             1 0 249
> mean(modtrain==dfttrain$Promoffer)
[1] 1
> mean(modtrain!=dfttrain$Promoffer)
[1] 0
>

```

Predict let us look at the classification matrix, now here you would see few errors.

(Refer Slide Time: 10:21)

```

444 cneck.names = F)
228
229 modtrain=predict(mod1, dfttrain[,c(3)], type = "class")
230 table("Actual value"=dfttrain$Promoffer, "Predicted value"=modtrain)
231 #classification accuracy
232 mean(modtrain==dfttrain$Promoffer)
233 #misclassification error
234 mean(modtrain!=dfttrain$Promoffer)
235
236 modvalid=predict(mod1, dfvalid[,c(3)], type = "class")
237 table("Actual value"=dfvalid$Promoffer, "Predicted value"=modvalid)
238 #classification accuracy
239 mean(modvalid==dfvalid$Promoffer)
240 #misclassification error
241 mean(modvalid!=dfvalid$Promoffer)
242
243 modltest=predict(mod1, dfitest[,c(3)], type = "class")

```

```

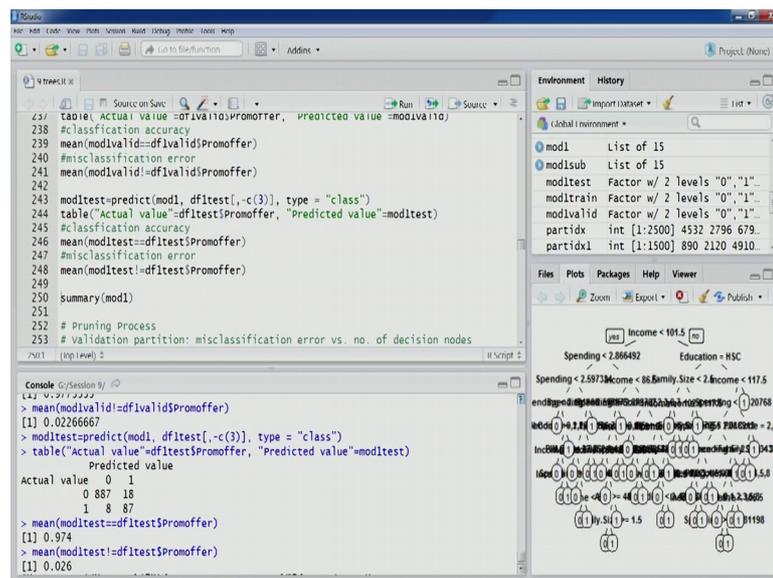
1 0 249
> mean(modtrain==dfttrain$Promoffer)
[1] 1
> mean(modtrain!=dfttrain$Promoffer)
[1] 0
> modvalid=predict(mod1, dfvalid[,c(3)], type = "class")
> table("Actual value"=dfvalid$Promoffer, "Predicted value"=modvalid)
      Predicted value
Actual value 0 1
             0 1342 22
             1 12 124
>

```

Right 22 class 0 records have been incorrectly classified as class 1 and 12 class 1 records have been incorrectly classified as class 0. So, let us look at the accuracy and misclassification error number. So, you can see 0.97. So, 97.70 percentage, that is; the accuracy and then we have 2.27 percentage that is the error right. So, So, when we apply the model.

So, there is this much difference in the performance of the model. So, it is of course, this is expected; because it is not possible for the; the over fitted 100 percent over fitted model to perform well on new data. However, the 97 this performance is also quite good; however, this is affected that it would be lower than the training partition. Now, let us look at the test partition and the performance. Now, let us look at the accuracy and misclassification error numbers we can see 97.4 in this case.

(Refer Slide Time: 11:29)



So, you if we compare from the previous performance; so would see then that in the validation partition is 97.7 and here it is that decreased 797.4. So, if we compare from the training to validation and test; so 100 to 97.7 to 97.4; so the performances though the performance is good for all the partition, but it is decreasing.

If we are interested in looking at the information; the way R output is presented. So, we can use the summary function as we have been doing for other techniques as well we can look at the details here. So, this is quite big output region being that, there are too many node numbers, because this is full grown tree. So, you can see in the summary output.

(Refer Slide Time: 12:25)

```
244 table(ACTUAL VALUE =dfitest$PROMOTTER, PREDICTED VALUE =moditest)
245 #classification accuracy
246 mean(modtest==dfitest$Promoffer)
247 #misclassification error
248 mean(modtest!=dfitest$Promoffer)
249
250 summary(mod1)
251
252 # Pruning Process
253 # Validation partition: misclassification error vs. no. of decision nodes
254 # Total no. of nodes in full grown tree
255 nrow(mod1$frame)
256 # No. of Decision Nodes
257 nrow(mod1$plits)
258 # No. of Terminal Nodes
259 nrow(mod1$frame)-nrow(mod1$plits)
260 # Node numbers
```

```
cp nsplit rel error
1 0.267068273 0 1.00000000
2 0.132530120 2 0.46586345
3 0.072289357 3 0.33333333
```

The call of this particular function details the number of observation then this is the CP table.

(Refer Slide Time: 12:32)

```
244 table(ACTUAL VALUE =dfitest$PROMOTTER, PREDICTED VALUE =moditest)
245 #classification accuracy
246 mean(modtest==dfitest$Promoffer)
247 #misclassification error
248 mean(modtest!=dfitest$Promoffer)
249
250 summary(mod1)
251
252 # Pruning Process
253 # Validation partition: misclassification error vs. no. of decision nodes
254 # Total no. of nodes in full grown tree
255 nrow(mod1$frame)
256 # No. of Decision Nodes
257 nrow(mod1$plits)
258 # No. of Terminal Nodes
259 nrow(mod1$frame)-nrow(mod1$plits)
260 # Node numbers
```

```
cp nsplit rel error
1 0.267068273 0 1.00000000
2 0.132530120 2 0.46586345
3 0.072289357 3 0.33333333
4 0.064257028 4 0.26104418
5 0.020080321 5 0.19678715
6 0.008032129 7 0.15662651
7 0.005020080 10 0.13253012
8 0.004016064 17 0.09236948
9 0.002409639 27 0.05220884
10 0.002008032 33 0.03614458
11 0.000000000 50 0.00000000
```

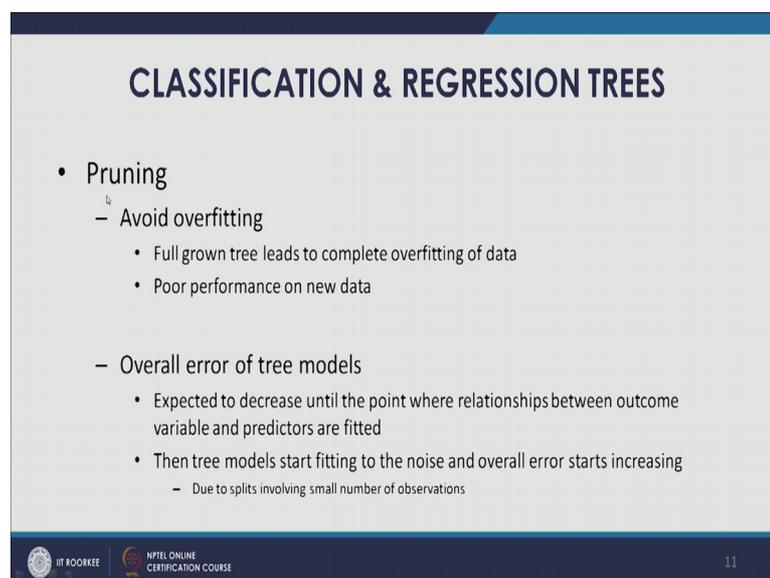
The CP table this is complexity a complexity parameter values. We will discuss this particular, what these values are about? In detail later on and split is the number of splits and then the relative error. So, I will discuss all these three all these three terms later on. Now let us look at the variable importance. So, for if we compare this to the previous data set that we had used the (Refer Time: 13:03) card data set we had just two variables

that, that is also there as you can see here you can see left son 19 1926 observation will go to the left son the right son the remaining 574 observation.

So, the same thing you can visualize in the tree diagram itself right. So, here; however, the node num node details this is not given here, because of the space problem in the tree diagram. Now, the primary split the split that has been performed is income less than 101.5 which we can see in the tree diagram as well the same information, how the improvement value that you can also see here after the split is performed. So, this is the improvement that is expected if this is split that is done right. Similarly, if you for node number two; we can look at other details 9 1926 observations and the complexity parameter value and then other details as we saw in the root node.

Now, look at this split spending less than 2.86 right improvement 3.31. So, node number 2. So, this is the node number 2. So, this split is performing in this fashion; using this predicted value combination you can look at the improvement and other things. So, similarly for all the node numbers in the summary output for all the node numbers, this detail we can find out right for all the node numbers this detail is available. So, further importance of some of these numbers that we have been talking about for example, complexity parameter and other things as we said we will discuss in the during the lecture itself. Now let us now that brings us to our next part that is pruning. So, let us go back to the slides.

(Refer Slide Time: 16:53)



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Avoid overfitting
 - Full grown tree leads to complete overfitting of data
 - Poor performance on new data
 - Overall error of tree models
 - Expected to decrease until the point where relationships between outcome variable and predictors are fitted
 - Then tree models start fitting to the noise and overall error starts increasing
 - Due to splits involving small number of observations

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE 11

So, as we have talked about we have discussed in the in some of the starting lectures of classification regression trees that there are two primary steps in cart procedures. In cart algorithm first one is recursive partitioning the second one is pruning. So, till now our discussion has revolved around the recursive partitioning, were we tried to create where we have tried to create the full grown tree we have tried to achieve your homogeneous subgroups right.

Now as we saw that this particular once we develop full grown tree it over fits the whole data and. So, how do we avoid this? So, pruning is the next step that can help us in avoiding over fitting. So, the idea is to avoid over fitting why we perform pruning? The idea is to avoid over fitting. So, full grown tree leads to complete over fitting of data as we have been saying poor performance on new data. So, because of the over fitting on new day new observation the performance will further decrease.

Now, if we if we look at the; the way tree models or in general other models that we have been talking about the focus is on typically the focus is on overall error. So, we always look to minimize overall error of classification models and also for tree models as well. So, what happens in when we try to you know minimize overall error? So, this particular error for any particular technique and also for tree models so, this is expected to decrease until the point where release relationships between outcome variables and predictors are fitted.

So, so, till the time the relationship between outcome variable and the set of predictors; that is being captured; that is being tapped in by the model. So, the till that time the overall error will continue to decrease right as we keep on you know building our model using that you know relationship using the information, that is there in those predictors. So, till the point the predictors information is being used to build the model the overall error will decrease; however, there is going to be a one point where the all the all that predictor information that could have been captured, that could be captured depending on the different strengths and weaknesses of different technique and for tree models.

As well at some point the models will start fitting to the noise once the all the predictors information; that could be useful for that classification tasks prediction task. In general, once that has been done the model will start fitting to the noise and overall error will start increasing, once we start fitting to the noise the overall error will start increasing. Now,

why specifically we talk about the tree models? Why that would happen? And in general also due to as you can see in the slide as well due to splits involving a small number of observations.

So, as we keep on doing splits using predictor value combinations. So, initially for every split will get you know much more reduction in overall error, but as we go along and most of the you know predictors information has been built into the model and once we start fitting to the noise we will also reach to the point where we are actually dealing with a small number of observation. So, we are further creating more and more splits most of the observation have already been put into their right rectangles right partitions. Now, further we are looking for few observations which are still to be correctly classified.

So, because those observations as we did through an exercise using (Refer Time: 20:51) cost data set the partition that we kept on creating using the graphic right. So, finally, we were also changing even if there is one point; that is, incorrectly classified in a particular group we also created a partition for even that one single point in that exercise right. So, you would see that essentially we are dealing with a small number of observation as we do to you know you know as we do too many splits right.

So, because of this the model tree models in general other models also they start fitting to the noise after some point once they predict information has been modelled. So, the overall error will start increase and will start fitting to the noise. So, the same thing happens in the when we develop full grown tree we keep on splitting we are we keep on splitting, we are keep on getting partition until we reach to the it will be generate pure homogeneous groups. So, some of the you know last splits some of the you know last split which have been created based on few observations. So, they are leading to they are actually fitting to the noise.

So, how do we overcome this how do we overcome this complete over fitting update of the data. So, there are two approaches. So, one is stop tree growth before it starts over fitting data or fitting noise all right.

(Refer Slide Time: 22:13)

CLASSIFICATION & REGRESSION TREES

- Pruning
 - Stop tree growth before it starts overfitting data or fitting noise
 - No. of splits or tree depth level
 - No. of observations in a node to attempt the split
 - Accepted level of reduction in impurity
 - Difficulties in determining the stopping point for such rules
 - Prune the full grown tree back to a level where it doesn't overfit data or fit noise
 - Use validation partition to prune the tree modeled with training partition
 - Idea is to remove the tree branches which don't reduce the error rate further

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE 12

So, we have to find out you know you know that point where the tree has stopped you know building the predictor. You stop using the predictors information to build the model, rather it has started the model as the technique has started over fitting the data or noise. So, how that can be done? So, number of is number of splits or tree depth level. So, if we can find out that optimum number of splits or tree depth level, where this actually can happen where this is actually going to happen the over fitting of data are fitting to the noise you know the number of split after which this starts to happen, if we are able to find out that number.

Then, probably we can stop tree growth there and we will get the optimized tree second approach could be number of observations in a node to attempt this split. So, how many observations could be there before and split is attempted. So, we are able to understand for a particular problem and the data set keeping everything in mind, if we are able to determine that and then also we can stop tree growth accepted a level of reduction in impurity. So, once we create split we are expecting certain reduction in the impurity the different two important impurity matrix we talked about gini index entropy measure. So, what is the level of reduction that we expect once we create this one split once we create partition and different subgroups.

So, if we are able to determine a; you know accepted level of you know impurity reaction in impurity. So, that can also be used to stop tree growth. So, this is another

approach now these three things that we talked about, whether it is number of a splits that can be performed optimized number of splits or optimal number of observation in a node to attend the split or the accepted level of reduction in impurity if. So, you would see with all these three you know approaches it is difficult to determine these points and how do we determine the number of splits or number of observation.

It is slightly difficult more number of experimentations would have to be performed right, it is going to take lot of time to build the model and to find the optimum you know these points and optimum values of some of these approaches. So, because of this difficulty this particular approach is you know not that popular and there is another approach; that is, generally tried which is actually the pruning approach that is the focus of this discussion this particular exercise also.

So, prune the second approach is prune the full grown tree back to a level where it does not over fit or fit noise. So, that is the second approach. So, once first we can develop the full grown tree and then we start pruning it back to the level where it does not over fit data or fit noise. So, that is one approach how it can be performed we can use validation partition. So, we can you know build our tree we can construct our tree using the training partition and, then we can start pruning it back using the validation partition.

So, this is one particular technique. So, we talked about in the starting you know lectures of this course, we talked about the importance of partitioning the training partition validation partition and also the test partition, but if we look at the some of the techniques that we have tried we have just been creating two partitions training partition and test partition.

Now, in this particular technique you would see the importance of validation partition training partition is being used to build the tree that is full grown tree validation partition is being used to define, the model to fine tune the model how to prune back the tree to the level where it is stops fitting to the noise or stops over fitting to the data. So, validation partition also is going to be in a way part of the modelling process. And therefore, the importance of test partition which would act as a new data partition and where we can apply our model and check it is performance.

So, in this particular technique you we see that the training partition and the validation partition they both become part of the modelling process. So, in this approach the

pruning approach; that we just discussed the idea is to remove the tree branches. When we talk about removing the; you know pruning that we back to a certain level. So, the idea is to remove the tree branches which do not reduce the error rate further right. So, the error rate is not you know some of the branches which are not able to decrease the error the overall error then probably we should remove those branches.

So, that is the main idea. So, so, to implement this particular approaches we need to we need to be able to understand to identify or determine the branches or the nodes which are not in decreasing the overall error further and then remove them off slip them off so with this so, this so, a few other things in pruning. So, as we talked about.

(Refer Slide Time: 27:54)

CLASSIFICATION & REGRESSION TREES

- Pruning
 - Prune the full grown tree back to a level where it doesn't overfit data or fit noise
 - Find the point where error rate on validation partition starts to increase
 - Cost complexity parameter or complexity parameter (CP) in CART algorithm
$$CP = Err + PF * TL$$
Where Err is misclassification error, PF is penalty factor for tree length (TL)
 - Minimum error tree
 - Tree with minimum misclassification error on validation partition

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE 13

I find the point where error rate on validation partition starts to increase. So, we talked about the point where you know construct the full grown tree and use the using retaining partition and the prune impact using the validation partition.

So, first we need to find out the point where error rate on validations we can build the model using partition and then apply, then a score the validation partition and check the error rate and find out at which point the error rate starts to increase; that is, the point that is the level at which till which we have to prune the tree. So, another concept related to the; this pruning is cost complexity parameter or complexity parameter CP; that is that is generally used in cart algorithm.

So, the complexity parameter values are typically used to perform the pruning in many implementations of cart. So, what is a complexity parameter? So, if you look at the definition here this complexity parameter that is CP is error; that is nothing, but the overall error that we talk talking about the misclassification error. So, CP is error plus PF into TL. PF is the penalty factor and TL is the tree length or tree size.

So, for tree size we introduce a penalty factor and error plus this penalty factor. So, that will give us the CP value. So, you know if there are two particular tree models two candidate tree models and they both have the same misclassification error, but if particular you know they have both these tree models have different tree sizes different number of nodes.

So, it is the smaller tree model that would be selected that would be based on the complexity parameter, because there is going to be penalty on this on three sides. So, the model having bigger tree length or high figure tea tree size would have additional you know additional value additional value added to this error component and therefore, it is error the CP value will would be more than the tree having similar having same misclassification error, but of a smaller size.

So, (Refer Time: 30:24) parameter or cost complexity parameter in a sense can be used to control the size of the tree to control the length of the tree. So, a complexity parameter can also be used to prune back the tree and as we have talked about that validation partition; that can be also used error rate and we can also identify from, where it starts to increase or the second approach is use the compute the complexity values and complexity parameter values and use that to find out the optimum length of the tree now.

So, we will stop at this point. And we will continue our discussion on pruning in the next lecture.

Thank you.