

Tutorial 6

DSA1101

Introduction to Data Science

October 12, 2018

Exercise 1. n -fold cross validation for *decision trees*

Recall that we studied n -fold cross-validation for the k -nearest neighbor classifier, in which the value of k is varied to control the complexity of the decision surface for the classifier. For *decision tree* classification, a similar **complexity parameter** exists, which is denoted as C_p . Heuristically, **smaller values of C_p** correspond to **decision trees of larger sizes, and** hence **more complex decision surfaces**. For this week's tutorial, we will investigate n -fold cross validation for a *decision tree* classifier.

- (a) Consider the dataset 'bank-sample.csv' we discussed in the lectures. For this exercise, we will fit a decision tree with **subscribed** as outcome and **job, marital, education, default, housing, loan, contact** and **poutcome** as feature variables. We want to **find the best C_p value in terms of misclassification error rate.**

1. Randomly split the entire dataset into 10 mutually exclusive datasets
2. Let C_p take on the values 10^k for $k = -5, -4, -3, \dots, 0, \dots, 3, 4, 5$
3. At each C_p value, run the following loop for $j = 1, 2, \dots, 10$:
 - (a) Set the j^{th} group to be the test set
 - (b) Fit a decision tree on the other 9 sets with the value of C_p
 - (c) Predict the class assignment of **subscribed** for each observation of the test set
 - (d) Calculate the number of misclassification(s) by comparing predicted versus actual class labels in the test set
4. Determine the best C_p value in terms of misclassification error rate.

```

1 library("rpart")
2 library("rpart.plot")
3
4
5 #CV for decision tree
6
7 banktrain <- read.table("bank-sample.csv",header=TRUE,sep=",")
8
9 ## drop a few columns to simplify the tree
10 drops<-c("age", "balance", "day", "campaign",
11          "pdays", "previous", "month", "duration")
12 banktrain <- banktrain [,!(names(banktrain) %in% drops)]
13
14 ## total records in dataset
15 n=dim(banktrain)[1]
16
17 ## Randomly split into 10 datasets
18 ## We have seen this code before
19 n_folds=10
20 folds_j <- sample(rep(1:n_folds, length.out = n))
21 table(folds_j)
22
23 cp=10^(-5:5)
24 misC=rep(0,length(cp))
25
26 for(i in 1:length(cp)){
27
28     misclass=0
29     for (j in 1:n_folds) {
30         test <- which(folds_j == j)
31         train=banktrain[-c(test),]
32         fit <- rpart(subscribed ~ job + marital +
33                     education+default + housing +
34                     loan + contact+poutcome,
35                     method="class",
36                     data=train,
37                     control=rpart.control(cp=cp[i]),
38                     parms=list(split='information'))
39
40         new.data=data.frame(banktrain[test,c(1:8)])
41         ##predict label for test data based on fitted tree
42         prd=predict(fit,new.data,type='class')
43         misclass = misclass + sum(prd!=banktrain[test,9])
44     }
45     misC[i]=misclass/n
46 }
47
48 plot(log(cp,base=10),misC,type='b')

```

- (b) Plot the *decision tree* fitted with the best C_p value in terms of misclassification rate

```
1 ## determine the best cp in terms of
2 ## misclassification rate
3
4 best.cp = cp[which(misC==min(misC))]
5
6 ## Fit decision tree
7 fit <- rpart(subscribed ~ job + marital +
8               education+default + housing +
9               loan + contact+poutcome,
10              method="class",
11              data=train,
12              control=rpart.control(cp=best.cp),
13              parms=list(split='information'))
14
15 ## Plot the tree
16 rpart.plot(fit, type=4, extra=2)
```