# Email Classification Using Classification Trees

Jasmine Barrera, Joe Sanchez, Aditya Ranade

March 25, 2019

## 1 Introduction

These days people get a lot of spam emails (the emails which have advertisements, money making schemes, chain letters, etc.) which can be a serious problem for the email company. If the email company does not put a system in place to filter out the spam emails, it can lose its customers and that could have serious implications on the company. In this project, we employ the classification and regression tree (CART) approach to classify a given email as spam or non spam based on various attributes like percentage of the word "free" in the email, percentage of exclamation marks in the email, etc. on the spam base data set available on the UCI Machine Learning Repository. In the database, we have 4601 observations on 58 variables out of which 1813 were considered spam and 2788 were considered as non spam. The data is readily available to us, hence is a secondary dataset with a mix of continuous and categorical variable.

## 2 Method

The selection of the best tree is done through the cross validation error. The tree with the lowest cross validation error is considered the best tree. However, the best tree might have too many splits and might be very difficult to understand. Hence, we find the optimal tree such that the cross validation error is slightly higher than the optimal tree but the splits are comparatively less. The process to select the optimal tree based on testing data set with 10 cross fold validation and the one standard error rule as follows.

Initially we consider the full tree and the list of sub trees based on training data set incorporating 10 fold cross validation. The optimal subtree is selected as the smallest subtree which has the cross validation error less than one standard error greater than the minimum cross validated error. Then we pass the testing sample through the optimal tree and calculate the misclassification error rate. An ideal tree model should have a low misclassification error rate on the test sample.

## 3 Analysis and Results

### 3.1 Classifier 1

For classifier 1, the optimal tree is tree 10 with 11 splits and 12 terminal nodes and the corresponding misclassification rate of the optimal tree is 10.87%. The optimal tree has 53 false positive cases which means 53 emails out of 1535 test samples were predicted to be spam but were observed to be not spam in reality. The tree has 114 false negative cases which means 114 out of 1535 test samples were predicted to be not spam but were observed to be spam in reality. Hence the optimal tree for classifier 1 has 3.45% false positive error rate and 7.43% false negative error rate. Please note if we follow the 1 SE rule strictly, the optimal tree is tree 14 with 19 splits and 20 terminal nodes with misclassification rate of 9.7%, false positive and false negative error rates 3.65% and 6.05% respectively.

For the first classifier we plot a subtree of the optimal tree with 7 splits and 8 terminal nodes which can be seen in figure 1 in the figures section. This subtree has total misclassification error rate 12.18% with false positive and false negative error rates 3.45% and 8.73% respectively.

Some of the variables used in the construction of the optimal tree for the first classifier are cfexc, crllongest, crlaverage, cfdollar, wffree, wfyour, wfremove, crltotal, wfhp, wfmoney, wfall, wfhpl, wfgeorge, wf000, wfinternet, wflabs, wf857.

The cross validation estimates of errors and training errors of sequence of pruned trees against the trees' complexity for the first classifier can be found in figure 3 in the figures section.

## 3.2   Classifier 2

For classifier 2, the optimal tree is tree 10 with 12 splits and 13 terminal nodes and the corresponding misclassification rate of the optimal tree is 15.30%. The optimal tree has 14 false positive cases which means 14 emails out of 1535 test samples were predicted to be spam but were observed to be not spam in reality. The tree has 221 false negative cases which means 221 out of 1535 test samples were predicted to be not spam but were observed to be spam in reality. Hence the optimal tree for classifier 2 has 0.91% false positive error rate and 14.39% false negative error rate. Please note if we follow the 1 SE rule strictly, the optimal tree is tree 15 with 24 splits and 25 terminal nodes with misclassification rate of 13.16%, false positive and false negative error rates 1.04% and 12.12% respectively.

For the second classifier we plot a subtree of the optimal tree with 7 splits and 8 terminal nodes which can be seen in figure 2 in the figures section. This subtree has total misclassification error rate 17% with false positive and false negative error rates 0.78% and 16.22% respectively.

Some of the variables used in the construction of the optimal tree for the second classifier are wfremove, crllongest, wf000, wfyour, crltotal, cfexc, crlaverage, wfgeorge, wfyou, wfhp, wfhpl, wfall, wfmoney, wfaddresses, wf415, wf857, cfdollar, wfcredit, wf85, wforder, wffree, wfmake, wfmail, wfre, wfreceive.

The cross validation estimates of errors and training errors of sequence of pruned trees against the trees' complexity for the second classifier can be found in figure 4 in the figures section.

## 4   Conclusion

Before comparing the cross validation estimates of errors of sequence of pruned trees against the trees' complexity, it is worthy to note that for classifier 1 (figure 3 in Figures section), the crossvalidation error is slightly more than training error which is on expected lines. Howerver for classifier 2 (figure 4 in Figures section), the crossvalidation error is considerably higher than the training error, which indicates classifier 2 model is overfitted.

Comparing the cross validation estimates of errors of sequence of pruned trees against the trees' complexity for the first and second classifier from figure 3 and 4, indicates the cross validation error for the first classifier has a smooth decreasing curve as compared to the cross validation error for the second classifier. Also, the overall curve for the first classifier is smoother than the second classifier. However, a huge decrease in the cross validation error is observed as for the second classifier as compared to the first classifier. The overall cross validation error is significantly lower for the first classifier as compared to the second classifier for trees with all the sizes.

Comparing the training errors of sequence of pruned trees against the trees' complexity for the first and second classifier from figure 3 and 4, it appears the training error for the first classifier has a big drop as compared to the training error for the second classifier. However, when we observe the scale, we can see the change is approximately same for both the classifiers.

If we go by just the misclassification error, we will use classifier 1 since it has lower misclassification error. However, if we are concerned with avoiding false positive cases, which means we want to avoid cases where a non spam email is classified as a spam email as that might cause us to miss out on important and authentic emails, we introduce a high penalty for the false positive case. This classifier then reduces the false positive cases but, in the process, the false negative cases increase. This also makes the model overfit. So, the overall misclassification error increases by some amount and hence the second classifier has overall misclassification error greater than the first classifier. Depending on the demands, we can we can choose classifier 1 or 2 accordingly.
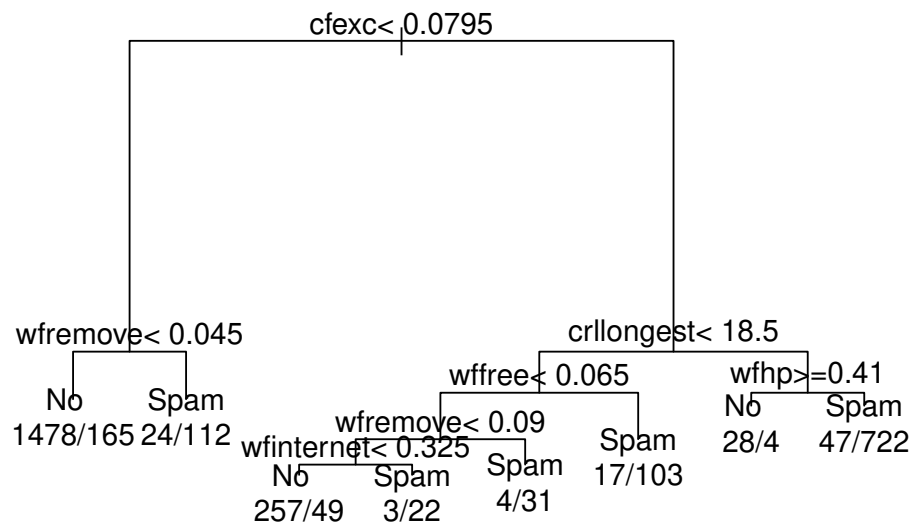
2

# 5 Figures

**Pruned Optimal Tree for Classifier 1**

cfexc< 0.0795

wfremove< 0.045

crllongest< 18.5

wffree< 0.065

wfhp>=0.41

No          Spam
1478/165  24/112

wfremove< 0.09

wfinternet< 0.325

Spam
17/103

No      Spam
28/4    47/722

No        Spam
257/49  3/22

Spam
4/31

Figure 1: Subtree of the optimal tree for first classifier

3

# Pruned Optimal Tree for Classifier 2

wfremove< 0.01

wf000< 0.335

wfgeorge>=0.08

cfexc< 0.6925

wfyou< 1.445

crllongest< 15.5

No
1776/369

No
9/28

Spam
0/123

No
14/0

wfhp>=0.15

No
39/44

Spam
1/131

No
7/5

Spam
12/508

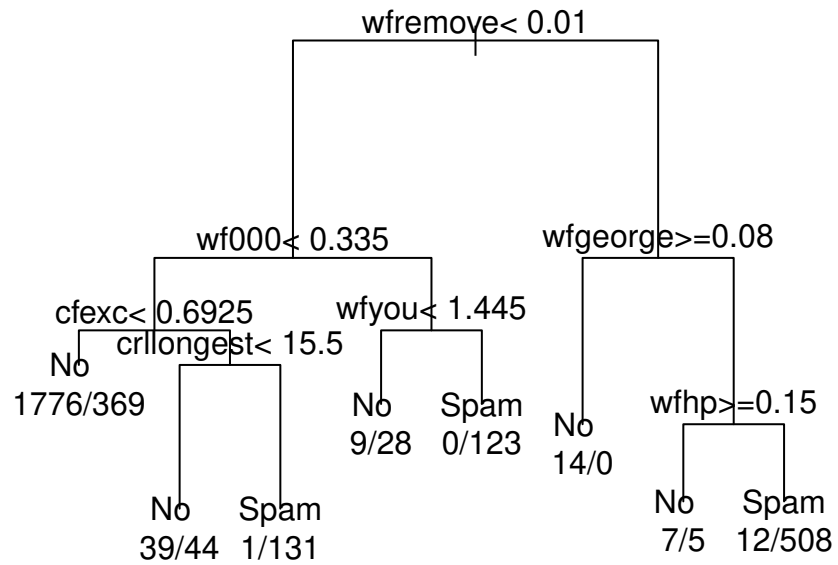Figure 2: Subtree of the optimal tree for second classifier

**Cross validation estimates of errors and training errors vs. tree complexity for classifier 1**
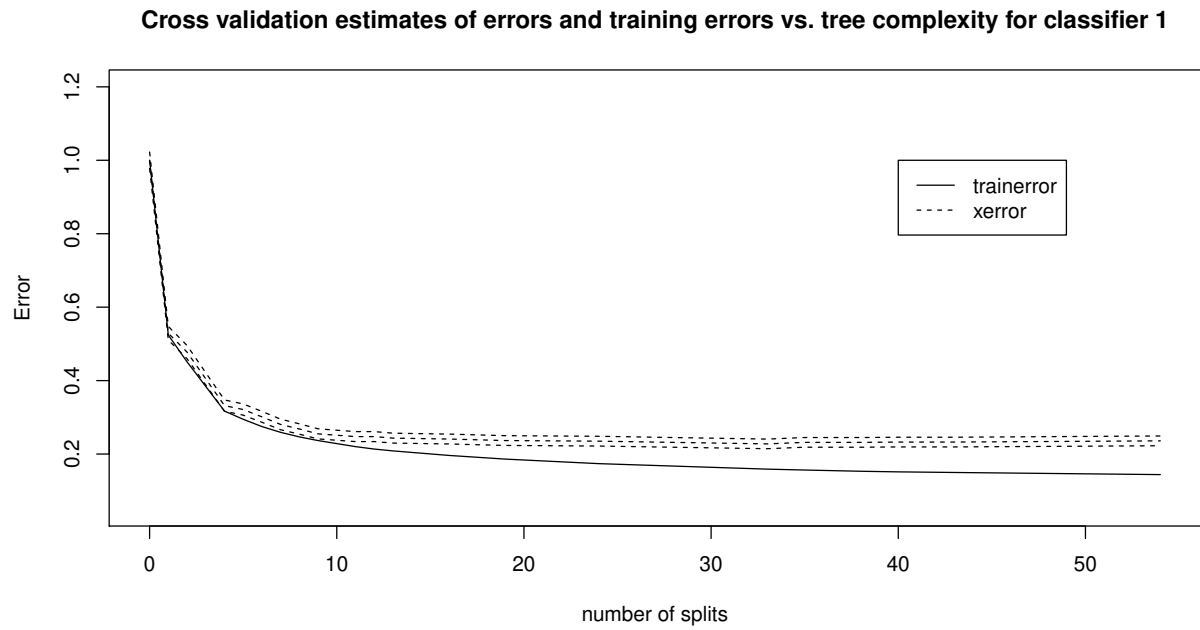


Figure 3: Crossvalidation estimates of errors and training errors for first classifier

# Appendix A: Contribution

Coding: Jasmine Barrera and Joe Sanchez
Report writing: Aditya Ranade

# Appendix B: R Code

```
1  library(rpart)
2  library(rpart.plot)
3  spam.data <- read.table("https://edoras.sdsu.edu/~jjfan/sta702/spamdata.txt",
4                          sep=",", na.strings="NA")
5  dim(spam.data)   # 4601    58
6  names(spam.data) <- c("wfmake", "wfaddress", "wfall", "wf3d", "wfour",
7                        "wfover", "wfremove", "wfinternet", "wforder", "wfmail",
8                        "wfreceive", "wfwill", "wfpeople", "wfreport", "wfaddresses",
9                        "wffree", "wfbusiness", "wfemail", "wfyou", "wfcredit", "wfyour",
10                       "wffont", "wf000", "wfmoney", "wfhp", "wfhpl", "wfgeorge", "wf650",
11                       "wflab", "wflabs", "wftelnet", "wf857", "wfdata", "wf415", "wf85",
12                       "wftechnology", "wf1999", "wfparts", "wfpm", "wfdirect", "wfcs",
13                       "wfmeeting", "wforiginal", "wfproject", "wfre", "wfedu", "wftable",
14                       "wfconference", "cfsc", "cfpar", "cfbrack", "cfexc", "cfdollar",
15                       "cfpound", "crlaverage", "crllongest", "crltotal", "spam")
16 table(spam.data$spam)   #0:2788 (nonspam); 1:1813 (spam) = data for spam column
17 spam.data$spam <- factor(spam.data$spam, levels=0:1,
18                          labels=c("No", "Spam"))
19
20
21 # split data into a training sample and a test sample, using stratified sampling
22 set1<-spam.data[spam.data$spam=="Spam",] #spam dataset
23 set0<-spam.data[spam.data$spam=="No",] #nonspam dataset
24 dim(set1)   # 1813    58
25 dim(set0)   # 2788    58
26 1813*2/3 #1208.667 emails = spam training
27 2788*2/3 #1858.667 emails = nonspam training
28
29 set.seed(858)
30 training1<-sample(1:1813,1208) #samples w/out replacement 1208 items from 1:1813
31 test1<-(1:1813)[-training1] #takes (1813*1/3) items from 1:1813
32 sum((1:1813)==sort(c(training1,test1))) # sort data to see if sum is same as 1:1813
33
34 training0<-sample(1:2788,1858)
35 test0<-(1:2788)[-training0]
36 sum((1:2788==sort(c(training0,test0)))) #2788
37
38
39 train<-rbind(set1[training1,],set0[training0,]) #combines training data for 0 and 1
40 test<-rbind(set1[test1,], set0[test0,])
41 dim(train)   # 3066    58
42 dim(test)    # 1535    58
43 3066+1535 #4601 = total
44 1208+1858 #3066 = training total
45 4601 - 3066 #1535 = testing total
46
47 ####### First Classifier #######
48 # tree growing and pruning with training data (aka fit classification tree to spam data
       using 10fold xvalidation)
49 my.control <- rpart.control(cp=0, xval=10) #cp = 0 means form largest tree; 10-fold xvalid
50 fit1<- rpart(spam ~ ., data=train, method="class",
51             control=my.control)  #this makes our largest tree and gives subtree sequence
52
```

```r
53  # plot the tree that corresponds to cp=0 (the largest, unpruned tree):
54  plot(fit1, margin=0.1)
55  text(fit1, use.n=T) #gives text to tree
56  title("Largest Tree for Classifier1")
57  printcp(fit1)  #recall, cp = cost complexity parameter
58  plotcp(fit1)
59
60  #tree 16 has absolutely smallest xerror (0.22765)
61  #0.22765 + 0.013098 (get 0.240748; look for xerror smaller than this)
62  #strictly following the 1SE rule, the optimal tree is tree 14
63  #but if we want simpler tree, we could go with tree 10.
64
65
66  #plot of xval estimates of error and training error against tree complexity
67  numsplits <- fit1$cptable[,2] #assigns the respective cptable values to variables names
68  trainerror <- fit1$cptable[,3]
69  xerror <- fit1$cptable[,4]
70  xstd <- fit1$cptable[,5]
71
72  plot(numsplits, trainerror, ylim=c(.05, 1.2), type="l") #plot training error (solid line)
73  lines(numsplits, xerror, lty =2) #plot xvalid error (dashed lines)
74  lines(numsplits, xerror-xstd, lty=2) #lower error bar for xvalid error
75  lines(numsplits, xerror+xstd, lty=2) #higher error bar for xvalid error
76  title("Cross-validation Error Estimates and Training Error for Classifier1")
77  legend(40, 1, c("trainerror","xerror"), lty=c(1,2)) #coordinates are for top right corner
        of legend
78
79  #get optimal tree by pruning (tree 10)
80  fit1pruned <- prune(fit1, cp=0.007)
81  print(fit1pruned)
82  plot(fit1pruned, margin=0.1)
83  text(fit1pruned, use.n=T)
84  title("Optimal Tree for Classifier1")
85  summary(fit1pruned)
86  summary(fit1, cp=0.007) #gives same result as above!
87
88  ##################misclassification error rate for tree 7##########################
89  #get subtree of optimal tree by pruning (tree 7)
90  fit1bpruned <- prune(fit1, cp=0.0135)
91  printcp(fit1bpruned)
92  plot(fit1bpruned, margin=0.1)
93  text(fit1bpruned, use.n=T)
94  title("Pruned Optimal Tree for Classifier1")
95  summary(fit1bpruned)
96  summary(fit1, cp=0.0135) #gives same result as above!
97
98  #running test data down pruned optimal tree (tree 7)
99  pred1<-predict(fit1bpruned,newdata=test,type="class")
100 error1<-table(test$spam,pred1)[1,2]+table(test$spam,pred1)[2,1]
101 error1 #total misclassification error
102 errorrate1<-error1/length(test$spam)
103 errorrate1 # total misclassification error rate=12.18%
104
105
106 #details of the misclassification error table (extra)
107 falserror1 <- table(test$spam,pred1) #table with number of false positive and false
        negative
108 falserror1 #930 non and 605 spam (test); 1011 non and 524 spam (pred1), where test and
        pred1 have total size 1535
109 falserrorate1 <- falserror1/length(test$spam)
110 falserrorate1 #divided false errors above by sample size 1535
111 #0.03452769 is the false positive error rate
112 #0.08729642 is the false negative error rate
113
114 ##################misclassification error rate for tree 10##########################
```

```r
115 #running test data down pruned optimal tree
116 pred1a<-predict(fit1pruned,newdata=test,type="class")
117 error1a<-table(test$spam,pred1a)[1,2]+table(test$spam,pred1a)[2,1]
118 error1a #total misclassification error
119 errorrate1a <-error1a/length(test$spam)
120 errorrate1a # total misclassification error rate=10.87%
121
122
123 #details of the misclassification error table (extra)
124 falserror1a <- table(test$spam,pred1a) #table with number of false positive and false
        negative
125 falserror1a #930 non and 605 spam (test); 877 non and 491 spam (pred1a), where test and
        pred1 have total size 1535
126 falserrorate1a <- falserror1a/length(test$spam)
127 falserrorate1a #divided false errors above by sample size 1535
128 #0.03452769 is the false positive error rate
129 #0.07426710 is the false negative error rate
130
131 #################misclassification error rate for tree 1#########################
132 fit1cpruned <- prune(fit1, cp=0.003)
133 #running test data down pruned optimal tree
134 pred1c<-predict(fit1cpruned,newdata=test,type="class")
135 error1c<-table(test$spam,pred1c)[1,2]+table(test$spam,pred1c)[2,1]
136 error1c #total misclassification error
137 errorrate1c <-error1c/length(test$spam)
138 errorrate1c # total misclassification error rate=9.70%
139
140
141 #details of the misclassification error table (extra)
142 falserror1c <- table(test$spam,pred1c) #table with number of false positive and false
        negative
143 falserror1c #930 non and 605 spam (test); 874 non and 512 spam (pred1a), where test and
        pred1 have total size 1535
144 falserrorate1c <- falserror1c/length(test$spam)
145 falserrorate1c #divided false errors above by sample size 1535
146 #0.03648208 is the false positive error rate
147 #0.06058632 is the false negative error rate
148
149
150
151 ####### Second Classifier #######
152 # tree growing and pruning with training data (aka fit classification tree to spam data
        using 10fold xvalidation)
153 my.control2 <- rpart.control(cp=0, xval=10) #cp = 0 means form largest tree; 10-fold xvalid
154 lmat <- matrix(c(0,10,1,0), byrow=T, nrow=2) #loss matrix, where we are penalize false
        positive 10 times more than false negative errors
155 fit2<- rpart(spam ~ ., data=train, method="class",
156              parm=list(loss=lmat),
157              control=my.control2)  #this makes our largest tree and gives subtree sequence
158
159 # plot the tree that corresponds to cp=0 (the largest, unpruned tree):
160 plot(fit2, margin=0.1)
161 title("Largest Tree for Classifier2")
162 text(fit2, use.n=T) #gives text to tree
163 printcp(fit2)  #recall, cp = cost complexity parameter
164 plotcp(fit2)
165 #tree 16 has absolutely smallest xerror (2.9007)
166 #2.9007 + 0.14490 (get 3.0456; look for xerror smaller than this)
167 #strictly following the 1SE rule, the optimal tree is tree 15; but if want smaller, could
        go with tree 10
168
169 #plot of xval estimates of error and training error against tree complexity
170 numsplits2 <- fit2$cptable[,2] #assigns the respective cptable values to variables names
171 trainerror2 <- fit2$cptable[,3]
172 xerror2 <- fit2$cptable[,4]
```

```
173 xstd2 <- fit2$cptable[,5]
174
175 plot(numsplits2, trainerror2, ylim=c(.05, 12), type="l") #plot training error (solid line)
176 lines(numsplits2, xerror2, lty =2) #plot xvalid error (dashed lines)
177 lines(numsplits2, xerror2-xstd2, lty=2) #lower error bar for xvalid error
178 lines(numsplits2, xerror2+xstd2, lty=2) #higher error bar for xvalid error
179 title("Cross-validation Error Estimates and Training Error for Classifier2")
180 legend(35, 10, c("trainerror2","xerror2"), lty=c(1,2)) #coordinates are for top right
        corner of legend
181
182 #get optimal tree by pruning (tree 10)
183 fit2pruned <- prune(fit2, cp=0.0135) #tree 10 with 12 splits and 13 terminal nodes
184 print(fit2pruned)
185 plot(fit2pruned, margin=0.1)
186 text(fit2pruned, use.n=T)
187 title("Optimal Tree for Classifier2")
188 summary(fit2pruned)
189 summary(fit2, cp=0.0135) #gives same result as above!
190
191 #get subtree of optimal tree by pruning (tree 6)
192 fit2bpruned <- prune(fit2, cp=0.035)
193 print(fit2bpruned)
194 plot(fit2bpruned, margin=0.1)
195 text(fit2bpruned, use.n=T)
196 title("Pruned Optimal Tree for Classifier2")
197 summary(fit2bpruned)
198 summary(fit2, cp=0.035) #gives same result as above!
199
200 #################misclassification error rate for tree 6#########################
201 #running test data down pruned optimal tree
202 pred2<-predict(fit2bpruned,newdata=test,type="class")
203 error2<-table(test$spam,pred2)[1,2]+table(test$spam,pred2)[2,1]
204 error2 #total misclassification error
205 errorrate2<-error2/length(test$spam)
206 errorrate2 #total misclassification error rate = 17%
207
208 #details of the misclassification error table (extra)
209 falserror2 <- table(test$spam,pred2) #table with number of false positive and false
        negative
210 falserror2 #930 non and 605 spam (test); 1167 non and 368 spam (pred1), where test and
        pred1 have total size 1535
211 falserrorate2 <- falserror2/length(test$spam)
212 falserrorate2 #divided false errors above by sample size 1535
213 #0.00781759 is the false positive error rate
214 #0.16221498 is the false negative error rate
215
216
217
218
219 #################misclassification error rate for tree 10#########################
220 #running test data down pruned optimal tree
221 pred2b<-predict(fit2pruned,newdata=test,type="class")
222 error2b<-table(test$spam,pred2b)[1,2]+table(test$spam,pred2b)[2,1]
223 error2b #total misclassification error
224 errorrate2b<-error2b/length(test$spam)
225 errorrate2b #total misclassification error rate
226
227 #details of the misclassification error table (extra)
228 falserror2b <- table(test$spam,pred2b) #table with number of false positive and false
        negative
229 falserror2b #930 non and 605 spam (test); 916 non and 384 spam (pred2b), where test and
        pred1 have total size 1535
230 falserrorate2b <- falserror2b/length(test$spam)
231 falserrorate2b #divided false errors above by sample size 1535
232 #0.00781759 is the false positive error rate
```

```
233  #0.16221498 is the false negative error rate
234
235  ###################misclassification error rate for tree 15##########################
236  fit2cpruned <- prune(fit2, cp=0.005)
237  printcp(fit2cpruned)
238  #running test data down pruned optimal tree
239  pred2c<-predict(fit2cpruned, newdata=test, type="class")
240  error2c<-table(test$spam, pred2c)[1,2]+table(test$spam, pred2c)[2,1]
241  error2c #total misclassification error
242  errorrate2c<-error2c/length(test$spam)
243  errorrate2c #total misclassification error rate = 13.16%
244
245  #details of the misclassification error table (extra)
246  falserror2c <- table(test$spam, pred2c) #table with number of false positive and false
            negative
247  falserror2c #930 non and 605 spam (test); 914 non and 419 spam (pred2c), where test and
            pred1 have total size 1535
248  falserrorate2c <- falserror2c/length(test$spam)
249  falserrorate2c #divided false errors above by sample size 1535
250  #0.01042345 is the false positive error rate
251  #0.12117264 is the false negative error rate
```

Listing 1: The R code used for analysis