

Task 2

Introduction

For task 2, we are going to deal with a dataset containing information about 4601 webmails. We have 48 variables describing the frequency of some specific words like “*remove*” in each observation, 6 variables describing the frequency of some specific chars like “\$” in one observation, and three variables, *capital_run_length_longest*, *capital_run_length_average* and *capital_run_length_total*, describing the length of the longest uninterrupted sequence of capital letters, the average length of uninterrupted sequences of capital letters, and the total number of capital letters in each observation respectively. We also have a variable called *spam*, which indicates whether this webmail is a spam with 0 and 1, where 1 for spam, and 0 for not spam. Here all our variables are numeric type. Our task is to use these 57 attribute variables to classify whether a webmail is spam.

Methodology

In order to validate the accuracy of our methods, we firstly divide our dataset into a train set, which contains 2500 observations, and a test set, which contains 2101 observations. We use the train set to train our models, and then apply it to the test set to validate its accuracy.

Results

Classification Trees

In this part, we are going to discuss the results obtained by complex tree model and pruned tree model. We begin with construct a complex tree model by dividing our observation into small non-overlapping regions according to some numerical criteria. Here we split our dataset until each leaf of our classification tree contains only less then 2 observations. The method used here is recursive binary splitting.

```
#grow complex tree using deviance as criterion
tree.mod = tree(factor(spam) ~ . ,data=train,
                control = tree.control(nobs = 2500, minsize = 2, mincut = 1),
                split = "deviance")

#tree.mod = rpart(spam ~ . ,data=train,
#                control = tree.control(nobs = 2500, minsize = 2, mincut = 1))

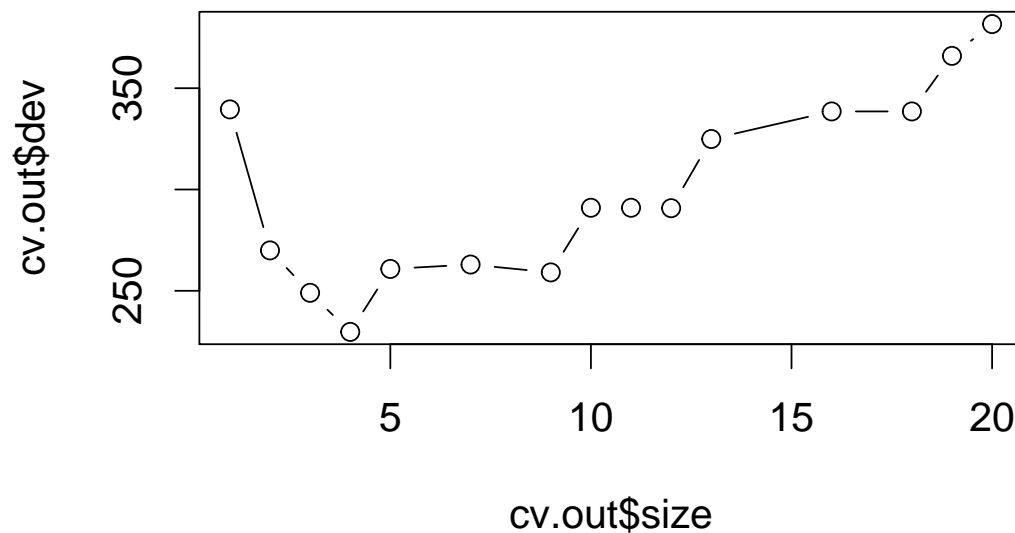
summary(tree.mod)
```

Classification tree:

```
tree(formula = factor(spam) ~ ., data = data.train, control = tree.control(nobs = 2500,
    minsize = 2, mincut = 1), split = "deviance")
```

Variables actually used in tree construction:


```
#use cross-validation to select tuning parameter for pruning the tree
set.seed(0829539)
cv.out=cv.tree(tree.mod,K=5)
par(cex=1.4)
plot(cv.out$size,cv.out$dev,type='b')
```



However, in this specific task, we can see that the cross-validation error is monotonously decreasing, so the optimal fold number is the original fold number. We choose best size=12 here such that the pruned tree model here is exactly the same as our complex tree model.

```
#prune the tree
prune.mod=prune.tree(tree.mod,best=12)
plot(prune.mod)
text(prune.mod,pretty=0)
```

Now we validate the accuracy of our classification tree model with the test data.

```
#make predictions on training and test set using the unpruned tree
pred.train<- predict(tree.mod,newdata=data.train)
classif.train<-ifelse(pred.train[,2]>=pred.train[,1],1,0)
err(data.train$spam,classif.train)

pred.test<-predict(tree.mod,newdata=data.test)
classif.test<-ifelse(1*pred.test[,2]>=pred.test[,1],1,0)
err(data.test$spam,classif.test)
```

```
#make predictions on training and test set using the pruned tree
pred.train<-predict(prune.mod,newdata=data.train)
classif.train<-ifelse(pred.train[,2]>=pred.train[,1],1,0)
err(data.train$spam,classif.train)
pred.test<-predict(prune.mod,newdata=data.test)
classif.test<-ifelse(1*pred.test[,2]>=pred.test[,1],1,0)
err(data.test$spam,classif.test)
```

We can conclude that our classification model performs very well. since the complex tree is the same as pruned tree here, we can see that the test error is just very slightly higher than the train error. There is not much overfitting here.

Different Classification Models

Now for this part, we are going to compare different classification models using two different classification scenarios.

Firstly, we consider the scenario where we have different prior probabilities and equal classification costs. We use the entire dataset to figure out the prior probability.

```
#(a)Account for different prior probabilities and equal classification costs
# (1) linear discriminant analysis
set.seed(0829539)
ldaS.out<-lda(spam~.,data=spamdata)
print(ldaS.out$prior)
```

```
      0      1
0.6059552 0.3940448
```

So, we take $P(\text{is spam})=0.394$, $P(\text{not spam})=0.606$ as our prior probability, and we can verify that their sum equals 1.

The four classification models we are going to compare is 1) Linear Discriminant Analysis, 2) Bagging, 3) Random Forests and 4) Gradient Boosting. Since the dimension here is very high and the mathematical assumption of our dataset is very weak, we do not suppose that Linear Discriminant Analysis without Principal Component analysis here will have very good performance. We just take it as a baseline. The three other methods are all methods used to improve the classification tree model, as tree models are always with high variance, thus high probability to cause overfitting.

The results are listed as below:

1) Linear Discriminant Analysis

```
lda.out<-lda(spam~.,prior=c(0.606,0.394),data=data.train)
```

Warning in lda.default(x, grouping, ...): variables are collinear

```
scaling <- lda.out$scaling %>% as.data.frame() %>% mutate(name = rownames(lda.out$scaling))
scaling %>% filter(name %in% c("char_freq_dollar", "word_freq_remove", "word_freq_table",
                             "word_freq_conference", "char_freq_dotcomma",
                             "char_freq_parenthesis", "char_freq_bracket", "char_freq_exclmark" ))
```

	LD1	name
1	0.6815563	word_freq_remove
2	-0.4296151	word_freq_table
3	-0.4725607	word_freq_conference
4	-0.5919005	char_freq_dotcomma
5	-0.4811911	char_freq_parenthesis
6	-1.3138634	char_freq_bracket
7	0.2859002	char_freq_exclmark
8	0.9359106	char_freq_dollar

Here we can see that the three attributes highly correlated with spam webmails—the frequency of “\$”, “!” and “remove”, are exactly the same as what we have got from our classification tree model. What is interesting here is we can also see that the frequency of “table”, “;” and “parenthesis” shows strong negative correlation with spam webmails.

```
pred.train<-predict(lda.out,data.train)
err(data.train$spam,pred.train$class)
```

```
pred.test<-predict(lda.out,data.test)
err(data.test$spam,pred.test$class)
```

2) Bagging

just some test

```
#(2) Bagging
set.seed(0829539)
bag.mod=randomForest(spam ~ ., data = data.train,
                      classwt=c(0.606,0.394),
                      mtry=57, ntree=nrotree, importance=TRUE)

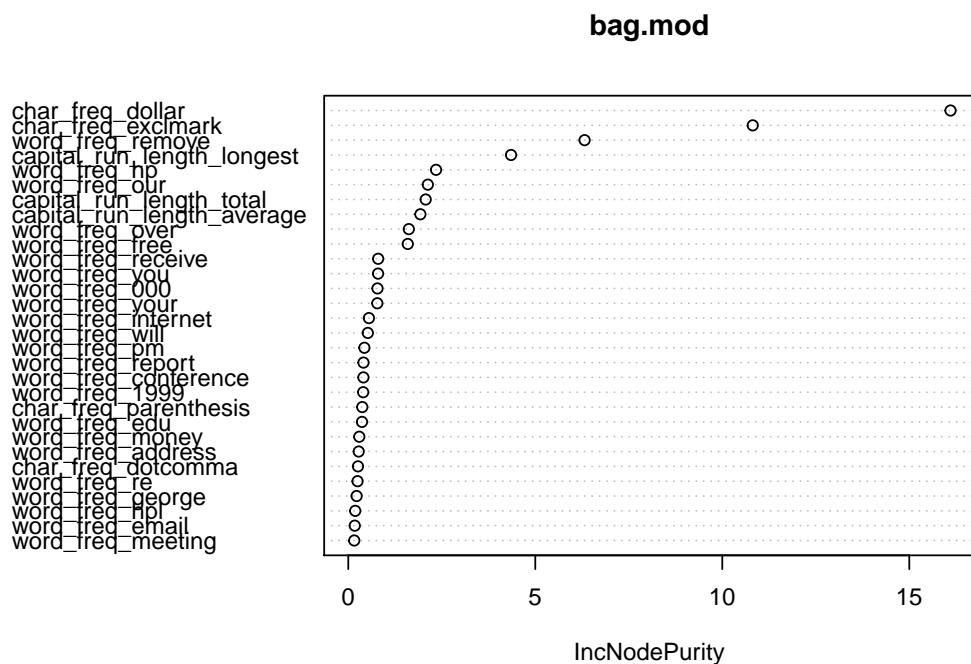
bag.mod
```

Call:

```
randomForest(formula = spam ~ ., data = data.train, classwt = c(0.606, 0.394), mtry = 57, ntree = 
              Type of random forest: regression
              Number of trees: 500
No. of variables tried at each split: 57
```

```
Mean of squared residuals: 0.07989729
% Var explained: 66.82
```

```
#plot variable importance
varImpPlot(bag.mod,type=2,cex=.8)
```



```
pred.train<-predict(bag.mod,newdata=data.train)
err(data.train$spam,pred.train)

pred.test<-predict(bag.mod,newdata=data.test)
err(data.test$spam,pred.test)
```

Actually, with the *MeanDecreaseGini* graph we can see that the frequency of “\$”, “!” and “remove” made evidently the most contribution. This is the same conclusion we have drawn from our classification tree model and Linear Discriminant Analysis. Here Out Of Bag error rate is 6.76%, which is acceptable.

3) Random Forests

just some text.

```
# (3) random forests
set.seed(0829539)
rf.mod=randomForest(spam~.,data=data.train,classwt=c(0.606,0.394),mtry=5,ntree=nrotree,importance=TRUE)
rf.mod
```

Call:

```
randomForest(formula = spam ~ ., data = data.train, classwt = c(0.606, 0.394), mtry = 5, ntree = 1000)
```

Type of random forest: regression

Number of trees: 500

No. of variables tried at each split: 5

Mean of squared residuals: 0.07646554

% Var explained: 68.24

```
pred.train<-predict(rf.mod,newdata=data.train)
err(data.train$spam,pred.train)
```

```
pred.test<-predict(rf.mod,newdata=data.test)
err_1 <- err(data.test$spam,pred.test)
```

```
kable(err_1$table)
```

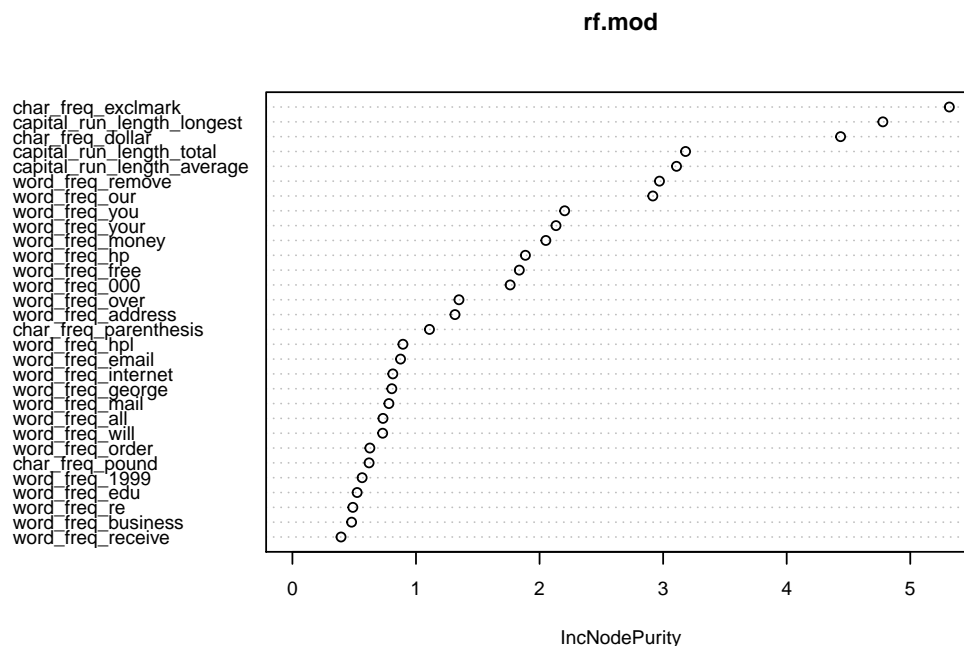
Warning in kable_pipe(x = structure(character(0), .Dim = c(0L, 0L), .Dimnames =
list(: The table should have a header (column names)

```
kable(err_1$numeric)
```

Warning in kable_pipe(x = structure(character(0), .Dim = c(0L, 0L), .Dimnames =
list(: The table should have a header (column names)

```
#variable importance plot
#two measures are reported:
#increase in MSE computed on OOB samples when removing variable
#decrease in node impurity (or increase in node purity) resulting from splits on the variable
importance(rf.mod)
```

```
varImpPlot(rf.mod,type=2,cex=.7)
```



Random Forests Model is also based on the idea of decorrelating trees. Here we choose $m=5$ for the size of our predictor set. can see that although there is slightly difference, “\$”, “!” and “remove” still stand for spam webmails, as well as the length of the longest uninterrupted sequence of capital letters.

4) Gradient Boosting

Different from Bagging and Random Forests, the number of trees for Gradient Boosting is not expected to be as large as possible, since using too many trees may cause overfitting for Gradient Boosting. We start with deciding the optimal number of trees by cross-validation, for the same reason explained before.

Note that since we are going to suppose Bernoulli distribution in our gbm function, we should firstly make sure that our variable for classification is numeric type.

```
# (4) gradient boosting
set.seed(0829539)
data.train$spam<-ifelse(data.train$spam=="1_yes",1,0)
data.test$spam<-ifelse(data.test$spam=="1_yes",1,0)
spamdata$spam<-ifelse(spamdata$spam=="1_yes",1,0)

#use distribution="bernoulli" for binary target
#interaction.depth=4, means that we fit a tree that uses 4 splits
#(and that includes at most a 4-th order interaction)
```