

Homework 5

Charles Hwang

Professor Matthews

STAT 388-001

12 November 2019

Exercise 1

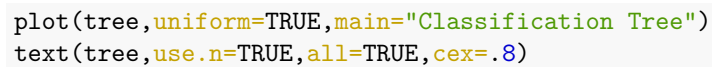
```
rm(list=ls())
train <- read.csv(file="/Users/newuser/Desktop/Notes/Undergraduate/STAT 338 - Predictive Analytics/titanic_train.csv")
test <- read.csv(file="/Users/newuser/Desktop/Notes/Undergraduate/STAT 338 - Predictive Analytics/titanic_test.csv")
train <- train[,-c(4,9,11)] # Removing categorical variables with large number of unique values ("Name", "Ticket")
test <- test[,-c(3,8,10)]
train$Sex <- as.factor(train$Sex)
test$Sex <- as.factor(test$Sex)
train$Embarked <- as.factor(train$Embarked)
test$Embarked <- as.factor(test$Embarked)
library(car)
library(gbm)
library(glmnet)
library(ISLR)
library(plotmo)
library(randomForest)
library(rpart)
set.seed(811)
```

Exercise 2

```
tree <- rpart(Survived~.,method="class",data=train)
printcp(tree)

##
## Classification tree:
## rpart(formula = Survived ~ ., data = train, method = "class")
##
## Variables actually used in tree construction:
## [1] Age      Embarked Fare      Pclass  Sex      SibSp
##
## Root node error: 342/891 = 0.38384
##
## n= 891
##
##          CP nsplit rel error  xerror    xstd
## 1 0.4444444      0  1.00000 1.00000 0.042446
## 2 0.030702      1  0.55556 0.55556 0.035750
```

```
plotcp(tree)
```



```

graph TD
    Root["Sex = U  
549/342"]
    Root -->|0| Node1["Age >= 6.5  
468/109"]
    Root -->|1| Node2["Pclass >= 2.5  
81/233"]
    Node1 -->|0| Node1L["460/93"]
    Node1 -->|1| Node1R["SibSp >= 2.5  
8/16"]
    Node1R -->|0| Node1RL["8/1"]
    Node1R -->|1| Node1RR["0/15"]
    Node2 -->|0| Node2L["Fare >= 23.35  
72/72"]
    Node2 -->|1| Node2R["9/161"]
    Node2L -->|0| Node2LL["24/3"]
    Node2L -->|1| Node2LR["Embarked = d  
48/69"]
    Node2LR -->|0| Node2LRL["Fare < 10.82  
32/31"]
    Node2LR -->|1| Node2LRR["16/38"]
    Node2LRL -->|0| Node2LRL0["22/15"]
    Node2LRL -->|1| Node2LRL1["Fare >= 17.6  
10/16"]
    Node2LRL1 -->|0| Node2LRL10["0"]
    Node2LRL1 -->|1| Node2LRL11["1"]
  
```

```
ptree <- prune(tree,cp=tree$cptable[which.min(tree$cptable[, "xerror"]), "CP"])
plot(ptree,uniform=TRUE,main="Pruned Classification Tree")
text(ptree,use.n=TRUE,all=TRUE)
```

Pruned Classification Tree



```
treepred <- predict(tree,test,type="class")
table(treepred)
```

```
## treepred
##  0  1
## 288 130
```

```
cat("Cross validation error:",min(tree$cptable[, "xerror"]))
```

```
## Cross validation error: 0.497076
```

Exercise 3

```
rf <- randomForest(as.factor(Survived)~.,data=train[!is.na(train$Age),],ntree=2000,importance=TRUE)
rf
```

```
##
```

```
## Call:
```

```
## randomForest(formula = as.factor(Survived) ~ ., data = train[!is.na(train$Age), ], ntree = 2000)
```

```
## Type of random forest: classification
```

```
## Number of trees: 2000
```

```
## No. of variables tried at each split: 2
```

```
##
```

```
## OOB estimate of error rate: 18.63%
```

```
## Confusion matrix:
```

```
## 0 1 class.error
```

```
## 0 382 42 0.0990566
```

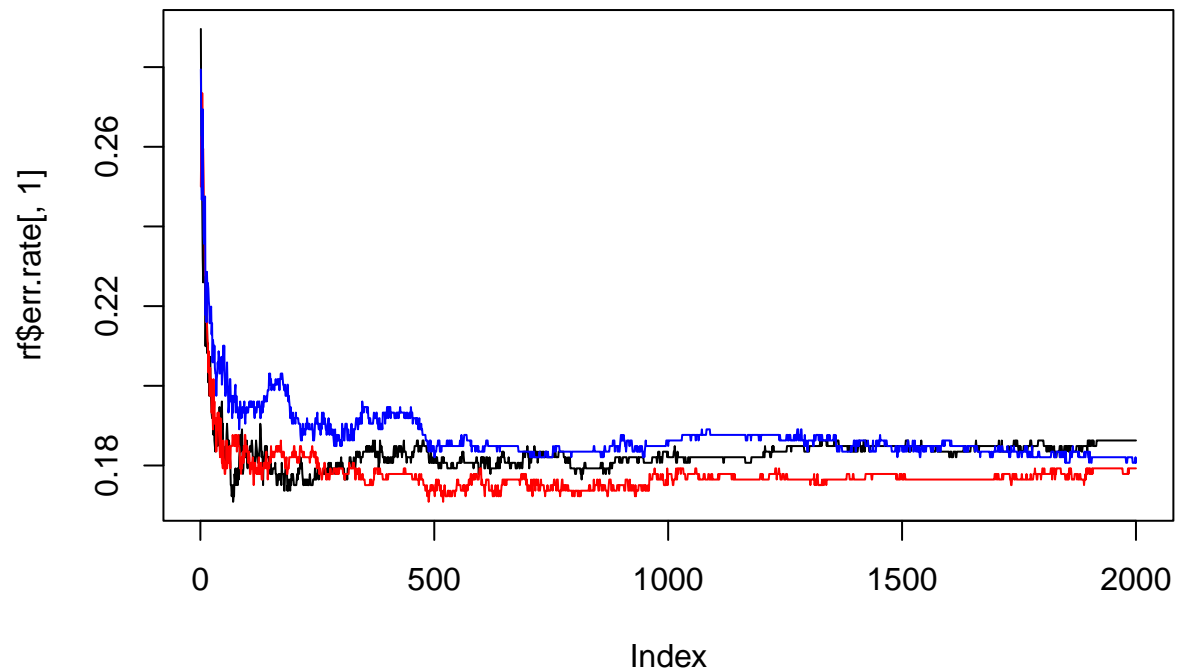
```
## 1 91 199 0.3137931
```

```
# The out-of-bag cross-validation error is 17.79%.
```

```
rf3 <- randomForest(as.factor(Survived)~.,data=train[!is.na(train$Age),],ntree=2000,importance=TRUE,mtry=3)
```

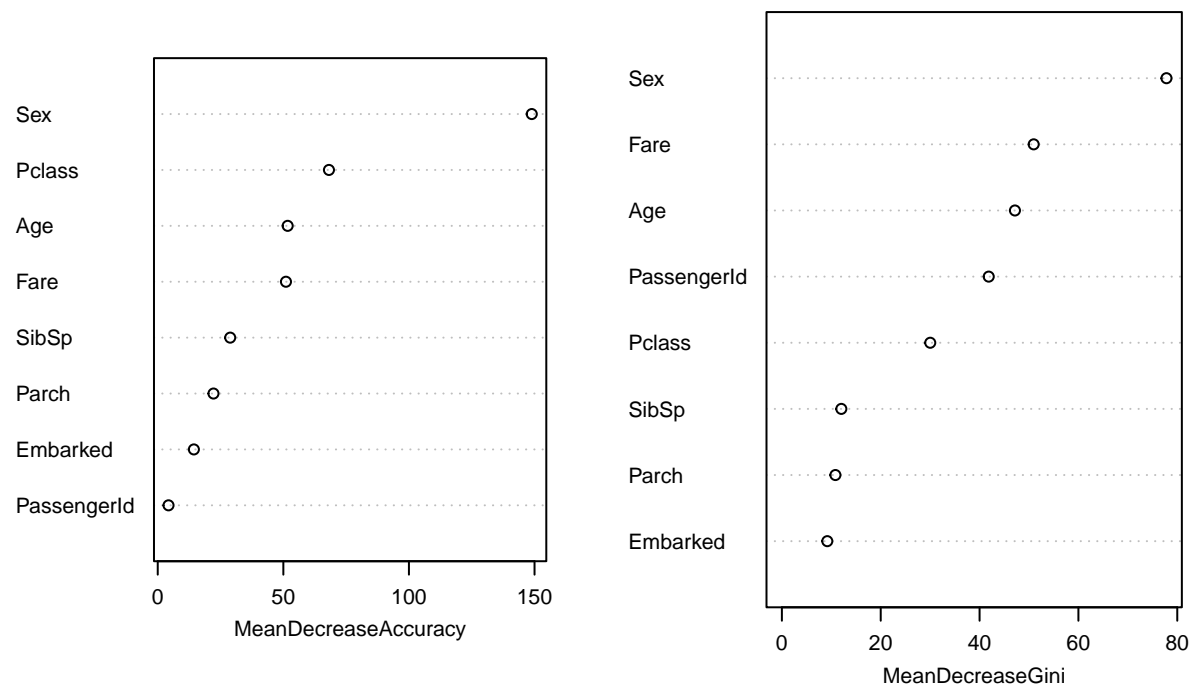
```
rf6 <- randomForest(as.factor(Survived)~.,data=train[!is.na(train$Age),],ntree=2000,importance=TRUE,mtry=6)
```

```
plot(rf$err.rate[,1],type="l")
points(1:2000,rf3$err.rate[,1],type="l",col="red")
points(1:2000,rf6$err.rate[,1],type="l",col="blue")
```



```
varImpPlot(rf,main="Variable Importance Plot",cex=.7)
```

Variable Importance Plot

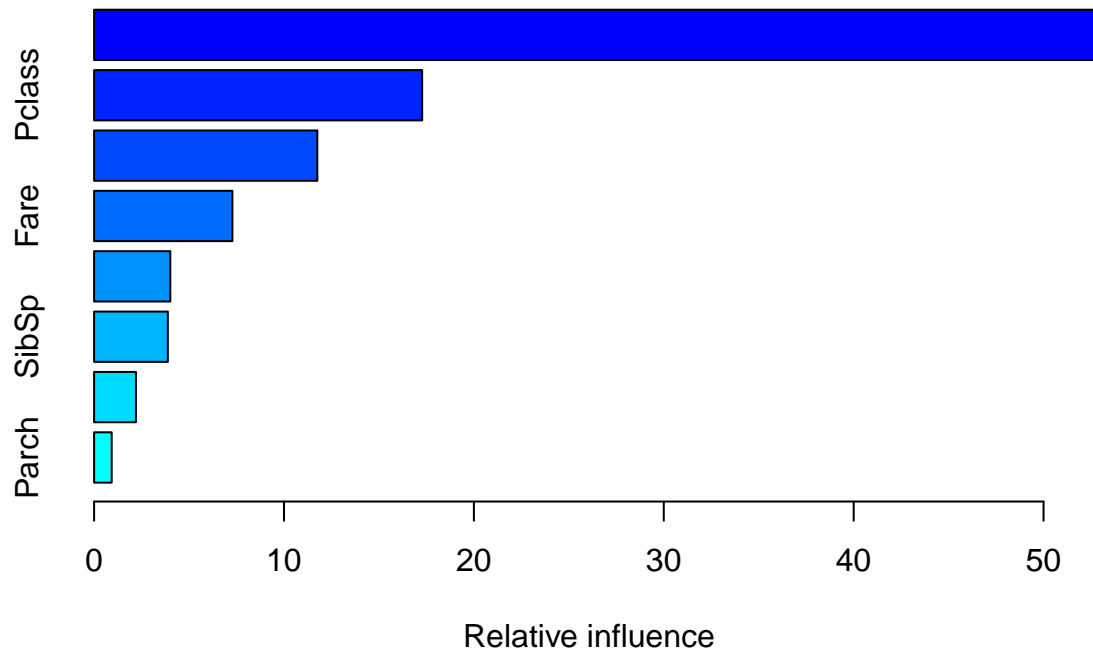


```
test <- rbind(train[1,-2],test) # See Stack Overflow for more information: https://stackoverflow.com/a/
test <- test[-1,]
rfpred <- predict(rf,test,type="class")
table(rfpred)
```

```
## rfpred
##    0    1
## 212 119
```

```
summary(gbm(Survived~.,data=train[!is.na(train$Age),]))
```

```
## Distribution not specified, assuming bernoulli ...
```



```
##          var    rel.inf
## Sex          Sex 52.6596638
## Pclass        Pclass 17.2761214
## Age           Age 11.7555255
## Fare          Fare  7.2836242
## PassengerId PassengerId 4.0143068
## SibSp         SibSp  3.8826045
## Embarked      Embarked  2.2028040
## Parch         Parch  0.9253498
```

```
rfpreddf <- as.data.frame(rep(0,length(rfpred))) # Exporting predictions to CSV file for Kaggle submission
rfpreddf$PassengerId <- test$PassengerId
rfpreddf$Survived <- rfpred
rfpreddf$`rep(0, length(rfpred))` <- NULL
Export(as.data.frame(rfpreddf),"Charles Hwang Submission.csv")
```

```
## Loading required namespace: rio
```

```
# Kaggle score: 290/418 (69.37799043% - nice!)
```