

INFX 573 Problem Set 7 - Prediction

Amol Surve

Due: Tuesday, November 29, 2015

Collaborators: Abhishek Gupta, Shreya Jain

Instructions:

Before beginning this assignment, please ensure you have access to R and RStudio.

1. Download the `problemset7.Rmd` file from Canvas. Open `problemset7.Rmd` in RStudio and supply your solutions to the assignment by editing `problemset7.Rmd`.
2. Replace the “Insert Your Name Here” text in the `author:` field with your own full name. Any collaborators must be listed on the top of your assignment.
3. Be sure to include well-documented (e.g. commented) code chunks, figures and clearly written text chunk explanations as necessary. Any figures should be clearly labeled and appropriately referenced within the text.
4. Collaboration on problem sets is acceptable, and even encouraged, but each student must turn in an individual write-up in his or her own words and his or her own work. The names of all collaborators must be listed on each assignment. Do not copy-and-paste from other students’ responses or code.
5. When you have completed the assignment and have **checked** that your code both runs in the Console and knits correctly when you click **Knit PDF**, rename the R Markdown file to `YourLastName_YourFirstName_ps7.Rmd`, knit a PDF and submit the PDF file on Canvas.

Setup:

In this problem set you will need, at minimum, the following R packages.

```
# Load standard libraries
library(tidyverse)
library(gridExtra)
library(MASS)
library(pROC)
library(arm)
library(randomForest)
library(xgboost)
```

Data: In this problem set we will use the `titanic` dataset used previously in class. The Titanic text file contains data about the survival of passengers aboard the Titanic. Table 1 contains a description of this data.

```
# Load data
titanicData <- read.csv('titanic.csv',header=T,na.strings=c(""))

sample<-titanicData
#retrieving the first 6 values
head(titanicData)
```

```
##      pclass survived                                name      sex
## 1         1         1                      Allen, Miss. Elisabeth Walton female
## 2         1         1                Allison, Master. Hudson Trevor    male
## 3         1         0                Allison, Miss. Helen Loraine female
## 4         1         0                Allison, Mr. Hudson Joshua Creighton  male
## 5         1         0 Allison, Mrs. Hudson J C (Bessie Waldo Daniels) female
## 6         1         1                      Anderson, Mr. Harry    male
##      age sibsp parch ticket      fare      cabin embarked boat body
## 1 29.0000      0      0 24160 211.3375      B5          S      2   NA
## 2  0.9167      1      2 113781 151.5500 C22 C26          S     11   NA
## 3  2.0000      1      2 113781 151.5500 C22 C26          S <NA>   NA
## 4 30.0000      1      2 113781 151.5500 C22 C26          S <NA>  135
## 5 25.0000      1      2 113781 151.5500 C22 C26          S <NA>   NA
## 6 48.0000      0      0 19952  26.5500      E12          S      3   NA
##      home.dest
## 1                      St Louis, MO
## 2 Montreal, PQ / Chesterville, ON
## 3 Montreal, PQ / Chesterville, ON
## 4 Montreal, PQ / Chesterville, ON
## 5 Montreal, PQ / Chesterville, ON
## 6                      New York, NY
```

```
#summarizing the data
summary(titanicData)
```

```
##      pclass      survived                                name
## Min.   :1.000   Min.   :0.000   Connolly, Miss. Kate      : 2
## 1st Qu.:2.000   1st Qu.:0.000   Kelly, Mr. James         : 2
## Median :3.000   Median :0.000   Abbing, Mr. Anthony      : 1
## Mean   :2.295   Mean   :0.382   Abbott, Master. Eugene Joseph : 1
## 3rd Qu.:3.000   3rd Qu.:1.000   Abbott, Mr. Rossmore Edward : 1
## Max.   :3.000   Max.   :1.000   Abbott, Mrs. Stanton (Rosa Hunt): 1
##                                     (Other)      :1301
##      sex      age      sibsp      parch
## female:466   Min.   : 0.1667   Min.   :0.0000   Min.   :0.000
## male :843    1st Qu.:21.0000   1st Qu.:0.0000   1st Qu.:0.000
##                                     Median :28.0000   Median :0.0000   Median :0.000
##                                     Mean   :29.8811   Mean   :0.4989   Mean   :0.385
##                                     3rd Qu.:39.0000   3rd Qu.:1.0000   3rd Qu.:0.000
##                                     Max.   :80.0000   Max.   :8.0000   Max.   :9.000
##                                     NA's   :263
##      ticket      fare      cabin      embarked
## CA. 2343: 11   Min.   : 0.000   C23 C25 C27      : 6   C :270
## 1601      : 8   1st Qu.: 7.896   B57 B59 B63 B66: 5   Q :123
## CA 2144 : 8   Median :14.454   G6              : 5   S :914
## 3101295 : 7   Mean   :33.295   B96 B98         : 4   NA's: 2
## 347077 : 7   3rd Qu.:31.275   C22 C26         : 4
## 347082 : 7   Max.   :512.329   (Other)         : 271
## (Other) :1261   NA's   :1      NA's         :1014
##      boat      body      home.dest
## 13      : 39   Min.   : 1.0   New York, NY      : 64
## C       : 38   1st Qu.: 72.0   London            : 14
## 15      : 37   Median :155.0   Montreal, PQ      : 10
## 14      : 33   Mean   :160.8   Cornwall / Akron, OH: 9
```

```
## 4      : 31  3rd Qu.:256.0  Paris, France      : 9
## (Other):308  Max.    :328.0  (Other)            :639
## NA's   :823  NA's    :1188  NA's                :564
```

```
#Converting "pclass" (Passenger Class Variable) to factors
titanicData$pclass <- as.factor(titanicData$pclass)
```

Variable	Description
pclass	Passenger Class (1 = 1st; 2 = 2nd; 3 = 3rd)
survived	Survival (0 = No; 1 = Yes)
name	Name
sex	Sex
age	Age
sibsp	Number of Siblings/Spouses Aboard
parch	Number of Parents/Children Aboard
ticket	Ticket Number
fare	Passenger Fare
cabin	Cabin
embarked	Port of Embarkation (C = Cherbourg; Q = Queenstown; S = Southampton)
boat	Lifeboat
body	Body Identification Number
home.dest	Home/Destination

Table 1: Description of variables in the Titanic Dataset

1. As part of this assignment we will evaluate the performance of a few different statistical learning methods. We will fit a particular statistical learning method on a set of *training* observations and measure its performance on a set of *test* observations.

- (a) Discuss the advantages of using a training/test split when evaluating statistical models.

Dividing the data into test and training set is the simplest partition possible for the cross sectional data which indicates that the available data fairly represents the real world processes we wish to model AND these real world processes are expected to remain relatively stable over time so that the model constructed using last months data would perform well for the the next month as well. The training data provides the raw material from which the predictive model is generated. Hence, we can see how well the model is performing on the available data and make changes iteratively in order to improve it's performance over the peroid of time based on the additional training data if available.

- (b) Split your data into a *training* and *test* set based on an 80-20 split, in other words, 80% of the observations will be in the training set.

```
dt = sort(sample(nrow(titanicData), nrow(titanicData)*.8))
#splitting the data into test and training set
#80% training data
#20% test data

#storing in dataframe
train<-titanicData[dt,]
#storing in dataframe
test<-titanicData[-dt,]

sapply(train,function(x) sum(is.na(x)))
```

```
##      pclass  survived      name      sex      age      sibsp      parch
##         0         0         0         0      208         0         0
##      ticket      fare      cabin embarked      boat      body home.dest
##         0         1        807         2      646      951      445
```

```
#observing the number of missing values in the dataframe
```

```
contrasts(titanicData$pclass)
```

```
##      2 3
## 1 0 0
## 2 1 0
## 3 0 1
```

```
train <- train[!is.na(train$pclass),]
rownames(train) <- NULL
```

2. In this problem set our goal is to predict the survival of passengers. First consider training a logistic regression model for survival that controls for the socioeconomic status of the passenger.

- (a) Fit the model described above using the `glm` function in R.

```
#Fitting the logistic regression model
results<-bayesglm(survived ~ pclass, family=binomial, data=train)

#Summary of the model
summary(results)
```

```
##
## Call:
## bayesglm(formula = survived ~ pclass, family = binomial, data = train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.3806  -0.7848  -0.7848   0.9870   1.6297
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   0.4660     0.1254   3.715 0.000203 ***
## pclass2      -0.7213     0.1832  -3.939 8.2e-05 ***
## pclass3      -1.4860     0.1577  -9.421 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1401.0  on 1046  degrees of freedom
## Residual deviance: 1304.6  on 1044  degrees of freedom
## AIC: 1310.6
##
## Number of Fisher Scoring iterations: 5
```

Note: I suggested `bayesglm` as well in case the model was unstable (you can see this with extremely large s.e. estimates for the coefficients). Be sure you included `pclass` as a `factor` because it is a categorical variable!

- (b) What might you conclude based on this model about the probability of survival for lower class passengers?

There is a positive coefficient estimate for the `pclass 1`. And `pclass2` and `pclass3` have negative coefficients which indicate that the passengers in lower classes have lesser chances of survival as compared to the higher class passengers. Hence, we can conclude that lower classes have a lower chance of survival as compared to the higher class.

3. Next, let's consider the performance of this model.

- (a) Predict the survival of passengers for each observation in your test set using the model fit in Problem 2. Save these predictions as `yhat`.

```
#predicting on the training data using the model
yhat<-predict(results, test, type="response")
```

- (b) Use a threshold of 0.5 to classify predictions. What is the number of false positives on the test data? Interpret this in your own words.

```
#classifying the predictions based on Above and Below 0.5 threshold value
yhat1<-ifelse(yhat>0.5,1,0)
#identifying the false positives
table(yhat1,test$survived)
```

```
##
## yhat1    0    1
##      0 150  55
##      1  21  36
```

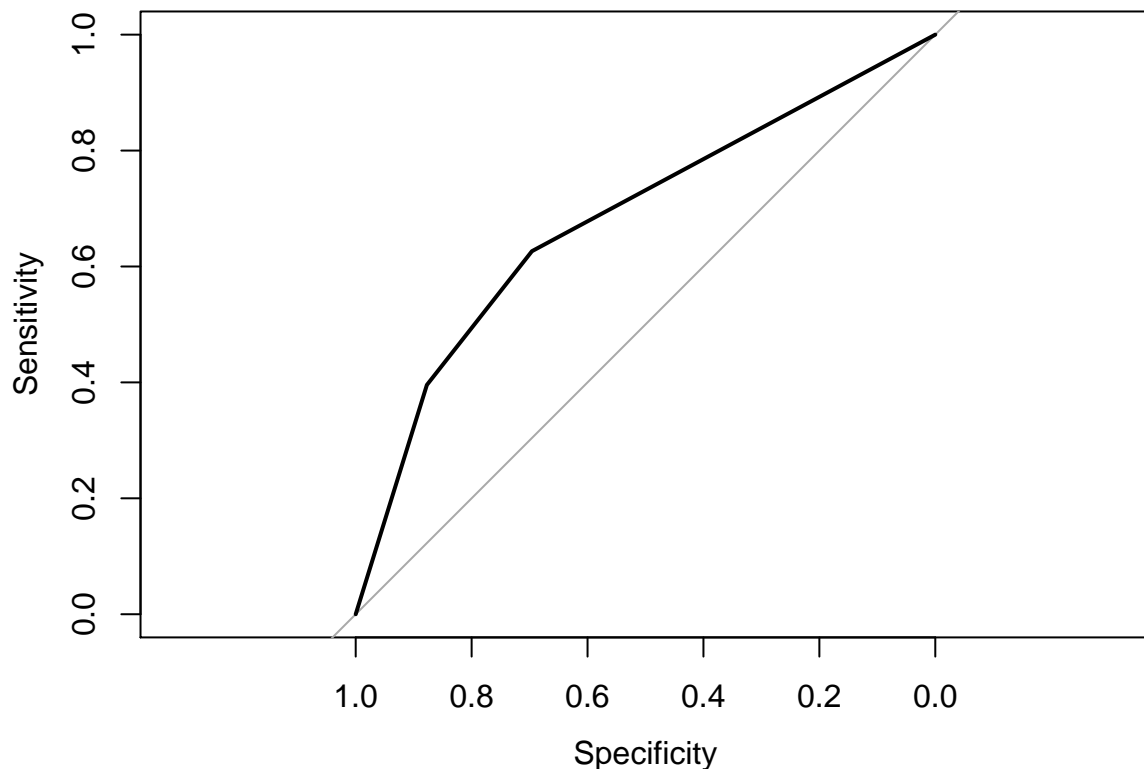
```
#identifying the percentage
mean(yhat1==test$survived)
```

```
## [1] 0.7099237
```

The number of false positives are 59. It is telling you that the person survived, even though he did not.

(c) Using the `roc` function, plot the ROC curve for this model. Discuss what you find.

```
#plotting the ROC curve using roc function
g <- roc(survived~yhat, data = test)
plot(g)
```



```
##
## Call:
## roc.formula(formula = survived ~ yhat, data = test)
##
## Data: yhat in 171 controls (survived 0) < 91 cases (survived 1).
## Area under the curve: 0.6828
```

The ROC curve demonstrates the following inferences: 1. It shows the trade-off between the sensitivity and specificity (any increase in sensitivity will be accompanied by decrease in specificity) 2. As you can see in the plot, the plot is closer to the sensitivity axis, which indicates better accuracy in the test results.

4. Suppose we use the data to construct a new predictor variable based on a passenger's listed title (i.e. Mr., Mrs., Miss., Master).

- (a) Why might this be an interesting variable to help predict passenger survival? Constructing a new predictor variable based on passenger's listed title i.e. Mr., Mrs., Miss., Master might provide us interesting insights about the survival of the passengers based on their age and sex. Following are some of the important questions that can be answered - 1. What proportion of the passengers were children women? 2. How many children and women were saved compared to all the people saved? 3. Were there any discrepancy between number of women and men saved? 4. What proportion of the Children were saved from each passenger class? Hence, these questions will provide us insights regarding the discrepancies in the number of passengers saved based on their age groups and sex.
- (b) Write a function to add this predictor to your dataset.

#the following code creates a new column

```
sample<-titanicData
sample$Title<-regmatches(as.character(sample$name),regexpr("\\.[A-z]{1,20}\\.", as.character(sample$name)))
sample$Title<-unlist(lapply(sample$Title,FUN=function(x) substr(x, 3, nchar(x)-1)))
table(sample$Title)
```

```
##
##          Capt          Col          Don          Dona          Dr
##           1           4           1           1           8
##   Jonkheer      Lady      Major      Master      Miss
##           1           1           2           61          260
##         Mlle         Mme          Mr          Mrs          Ms
##           2           1          757          197           2
##          Rev        Sir the Countess
##           8           1           1
```

```
sample$Title[which(sample$Title %in% c("Mme", "Mlle"))] <- "Miss"
sample$Title[which(sample$Title %in% c("Lady", "Ms", "the Countess", "Dona"))] <- "Mrs"
sample$Title[which(sample$Title=="Dr"&sample$sex=="female")] <- "Mrs"
sample$Title[which(sample$Title=="Dr"&sample$sex=="male")] <- "Mr"
sample$Title[which(sample$Title %in% c("Capt", "Col", "Don", "Jonkheer", "Major", "Rev", "Sir"))] <- "Mr"
sample$Title<-as.factor(sample$Title)
#convert to factor variable
str(sample)
```

```
## 'data.frame':   1309 obs. of  15 variables:
## $ pclass   : Factor w/ 3 levels "1","2","3": 1 1 1 1 1 1 1 1 1 ...
## $ survived : int   1 1 0 0 0 1 1 0 1 0 ...
## $ name      : Factor w/ 1307 levels "Abbing, Mr. Anthony",...: 22 24 25 26 27 31 46 47 51 55 ...
## $ sex       : Factor w/ 2 levels "female","male": 1 2 1 2 1 2 1 2 1 2 ...
## $ age       : num   29 0.917 2 30 25 ...
## $ sibsp     : int   0 1 1 1 1 0 1 0 2 0 ...
## $ parch     : int   0 2 2 2 2 0 0 0 0 0 ...
## $ ticket    : Factor w/ 929 levels "110152","110413",...: 188 50 50 50 50 125 93 16 77 826 ...
## $ fare      : num   211 152 152 152 152 ...
## $ cabin     : Factor w/ 186 levels "A10","A11","A14",...: 44 80 80 80 80 150 146 16 62 NA ...
## $ embarked  : Factor w/ 3 levels "C","Q","S": 3 3 3 3 3 3 3 3 1 ...
## $ boat      : Factor w/ 27 levels "1","10","11",...: 12 3 NA NA NA 13 2 NA 27 NA ...
## $ body      : int   NA NA NA 135 NA NA NA NA NA 22 ...
## $ home.dest : Factor w/ 369 levels "?Havana, Cuba",...: 309 231 231 231 231 237 162 24 22 22 ...
## $ Title     : Factor w/ 4 levels "Master","Miss",...: 2 1 2 3 4 3 2 3 4 3 ...
```

- (c) Fit a second logistic regression model including this new feature. Use the `summary` function to look at the model. Did this new feature improve the model?

```

dt1 = sort(sample(nrow(titanicData), nrow(titanicData)*.8))
#splitting the sample data into train and test
#storing in dataframe
trainNew<-sample[dt,]
#storing in dataframe
testNew<-sample[-dt,]
#fitting the logistic regression model
resultsNew<-glm(survived~pclass+Title, family=binomial, data=trainNew)
#summarizing the model
summary(resultsNew)

```

```

##
## Call:
## glm(formula = survived ~ pclass + Title, family = binomial, data = trainNew)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.2505  -0.6268  -0.4187   0.6369   2.2261
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.5637     0.3528   4.432 9.32e-06 ***
## pclass2        -0.9573     0.2311  -4.142 3.45e-05 ***
## pclass3        -1.8196     0.2051  -8.873 < 2e-16 ***
## TitleMiss       0.6258     0.3398   1.841  0.0656 .
## TitleMr        -2.1342     0.3305  -6.458 1.06e-10 ***
## TitleMrs       0.8858     0.3685   2.404  0.0162 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1401.0  on 1046  degrees of freedom
## Residual deviance:  956.4  on 1041  degrees of freedom
## AIC: 968.4
##
## Number of Fisher Scoring iterations: 4

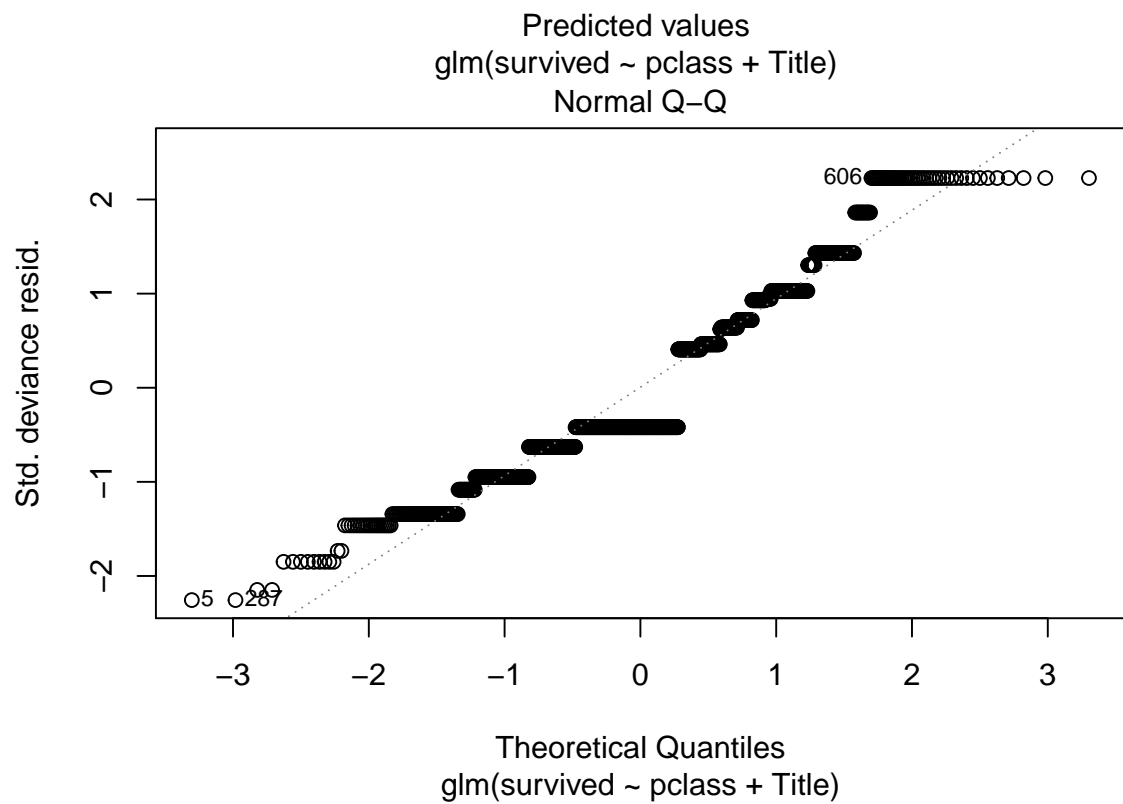
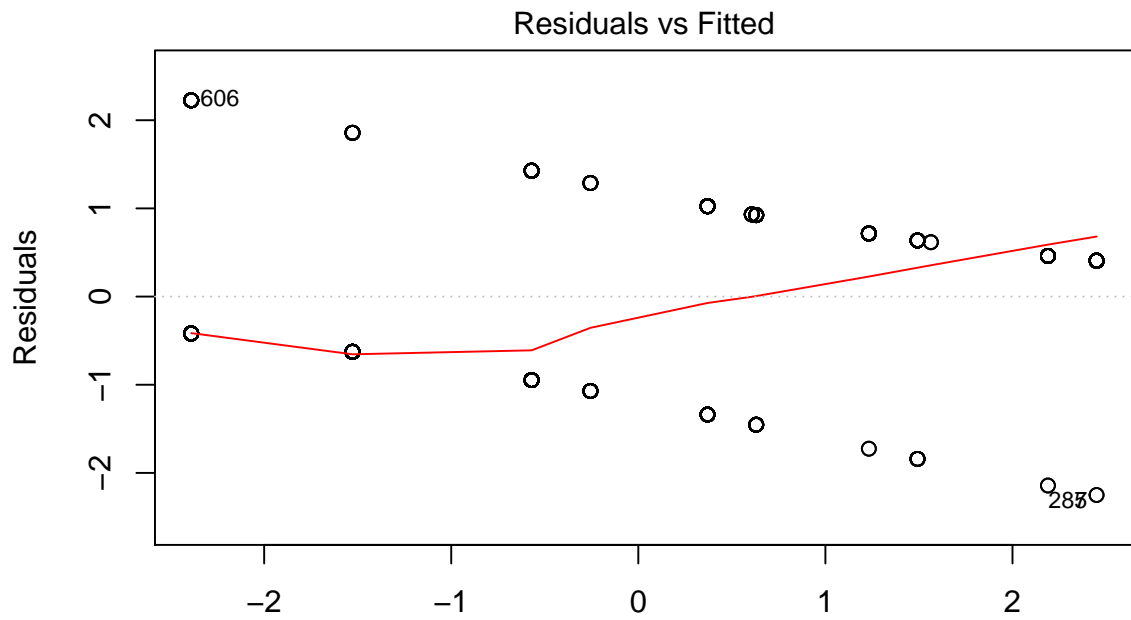
```

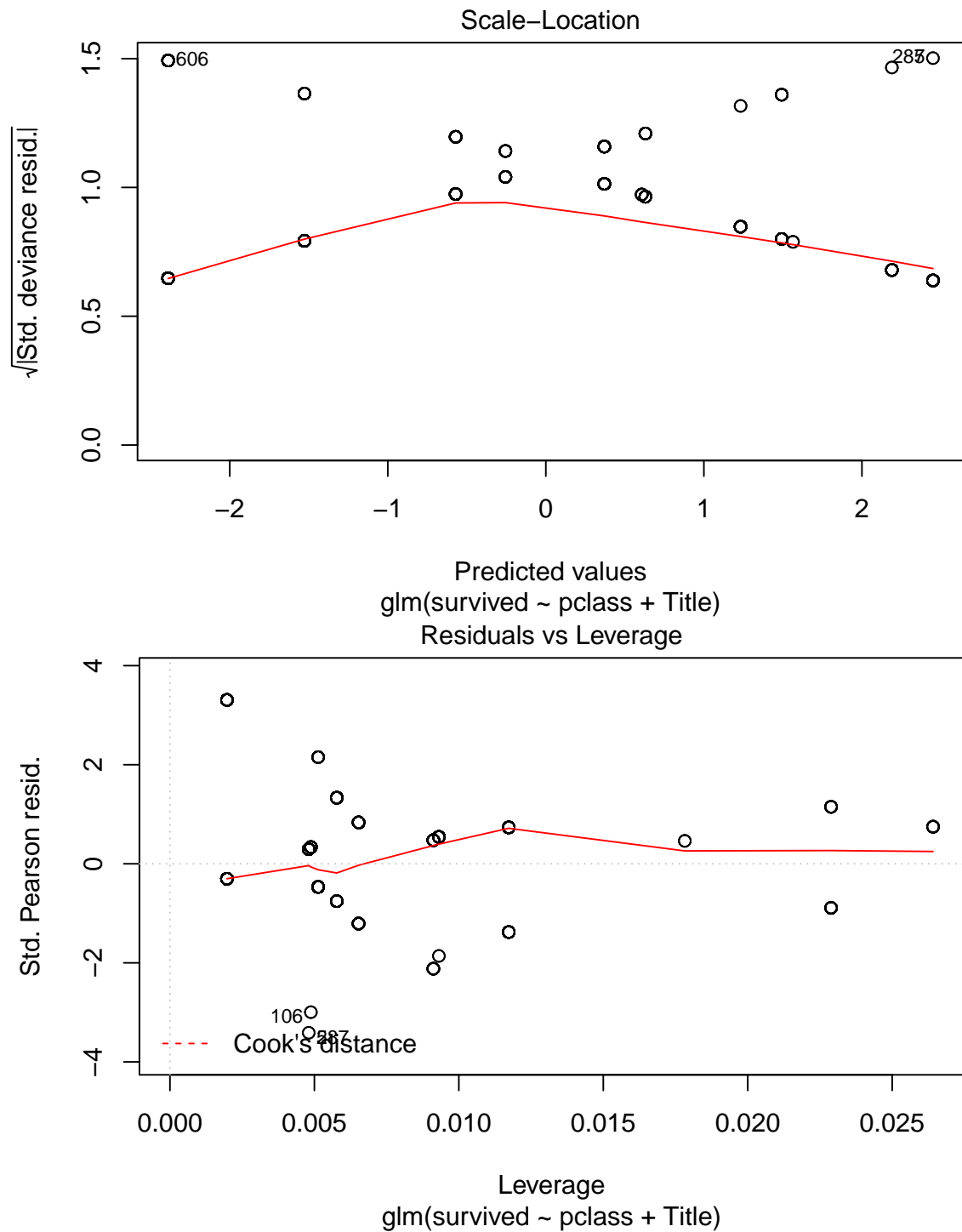
- (d) Comment on the overall fit of this model. For example, you might consider exploring when misclassification occurs.

```

plot(resultsNew)

```



#Residual Plots

The model fit is better compared to the previous model. There is some variance and also some of the high leverage points that are stretching the line. This model is better than the first model, but there is still missclassification based on the training datasets.

- (e) Predict the survival of passengers for each observation in your test data using the new model. Save these predictions as `yhat2`.

```
#predicting on the training data using the model
yhatlogi<-predict(resultsNew, testNew, type="response")
table(yhatlogi,testNew$survived)
```

```
##
## yhatlogi          0  1
## 0.0839378275295737 81 13
## 0.178331376963359 30  1
## 0.361135657303399 20 12
## 0.436388415681775  7  5
## 0.591451451496034 19 10
## 0.652490988517498 12  6
## 0.774217906196628  1 13
## 0.816426371508083  0  7
## 0.826888901206293  0  2
## 0.899306100011745  1  8
## 0.92053044816532  0 14
```

```
yhat11<-ifelse(yhat>0.5,1,0)
table(yhat11,testNew$survived)
```

```
##
## yhat11    0    1
##          0 150  55
##          1  21  36
```

```
mean(yhat11==testNew$survived)
```

```
## [1] 0.7099237
```

Percentage is 65.64 which is better than the previous model.

5. Another very popular classifier used in data science is called a *random forest*¹.

- (a) Use the `randomForest` function to fit a random forest model with passenger class and title as predictors. Make predictions for the test set using the random forest model. Save these predictions as `yhat3`.

```
#Fitting the Random forest Model
fit <- randomForest(as.factor(survived) ~ pclass+Title,
                    data=trainNew,
                    importance=TRUE,
                    ntree=1000)
#Displaying the model
fit
```

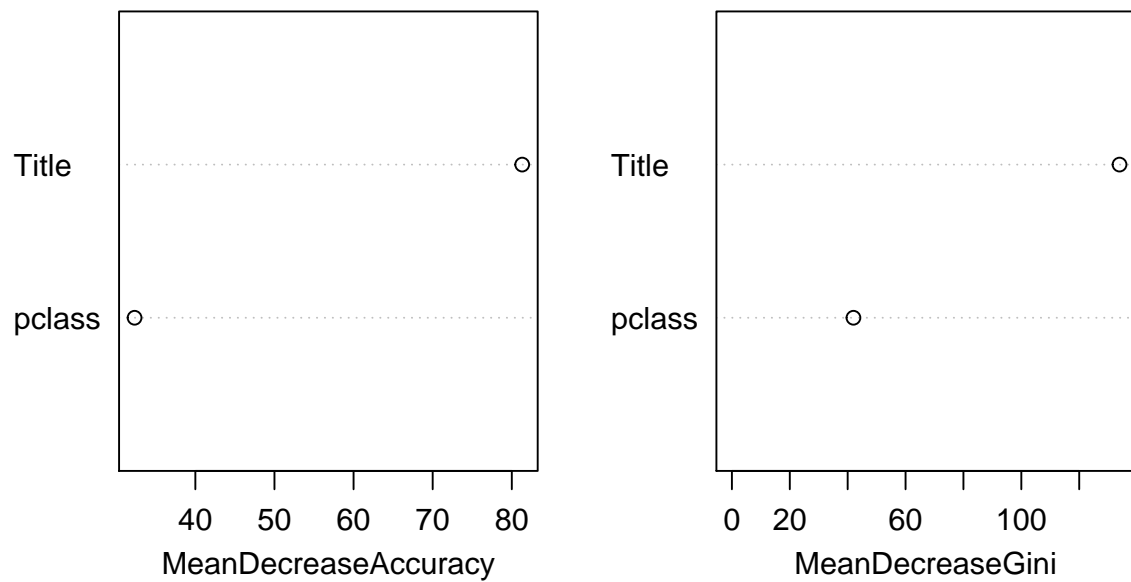
```
##
## Call:
## randomForest(formula = as.factor(survived) ~ pclass + Title,      data = trainNew, importance=TRUE,
##               Type of random forest: classification
##               Number of trees: 1000
##               No. of variables tried at each split: 1
##
##               OOB estimate of  error rate: 22.16%
```

¹https://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.htm

```
## Confusion matrix:
##      0   1 class.error
## 0 544  94  0.1473354
## 1 138 271  0.3374083
```

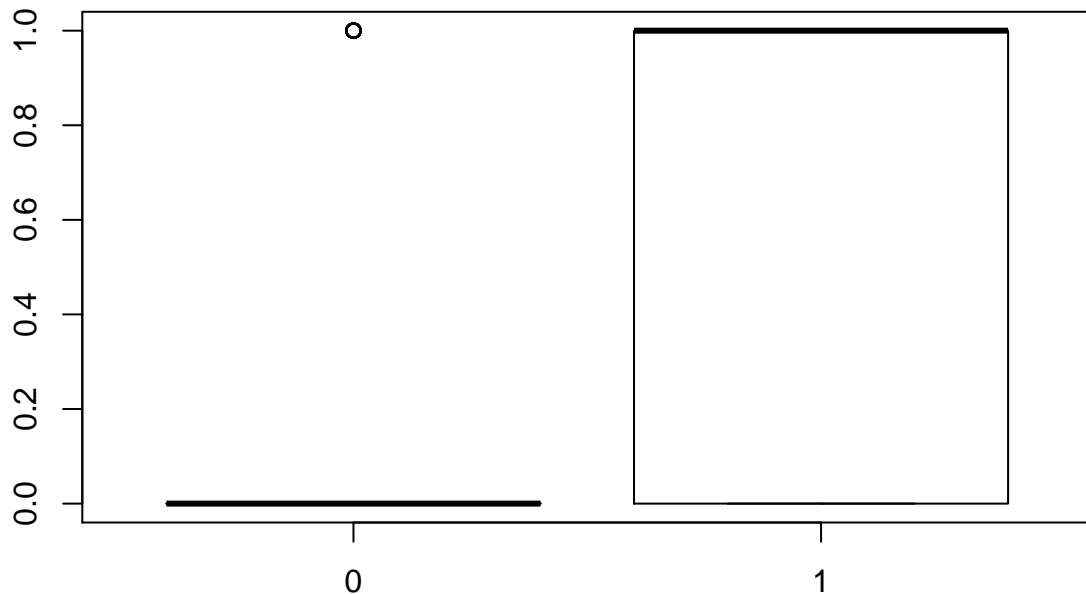
```
#plotting
varImpPlot(fit)
```

fit



```
#predicting the test set
yhatRandom<-predict(fit,newdata=testNew)
yhatRandomROC<-predict(fit,newdata=testNew,type="prob")[,2]

#plotting the Random Forest
plot(yhatRandom,testNew$survived)
```



```
#percentage of correct predictions
mean(yhatRandom==testNew$survived)
```

```
## [1] 0.7557252
```

Random forest Model is fitted using pclass and title as predictors. the accuracy of the model is 81.29%

- (b) Develop your own random forest model, attempting to improve the model performance. Make predictions for the test set using your new random forest model. Save these predictions as `yhat4`. In order to develop a new model, I figured out the variables with no missing values. also, variable like fare just had one missing value. So, I first removed these missing values from the training data and then fitted the random forest model on this training data. It helped to increase the accuracy of the previous model from 81.29% to 83.58%.

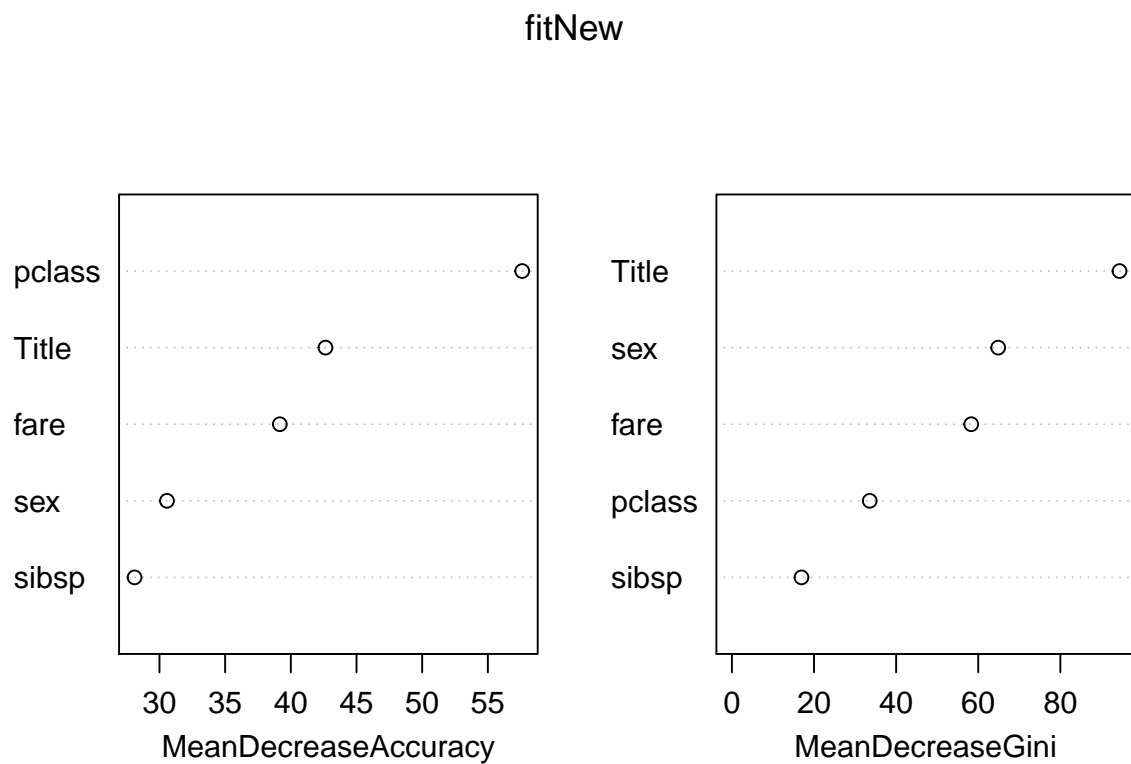
```
df<-trainNew
#copying into new dataframe
df<-df[!(is.na(df$fare) | df$fare==""), ]
#removeing na values from fare
sapply(df,function(x) sum(is.na(x)))
```

```
##      pclass  survived      name      sex      age      sibsp      parch
##         0         0         0         0      208         0         0
##      ticket      fare      cabin embarked      boat      body home.dest
##         0         0        806         2      645      951      444
##      Title
##         0
```

```
#checking whether fare has zero missing values
#Fitting the Random forest Model
fitNew <- randomForest(as.factor(survived) ~ pclass+Title+sex+sibsp+fare,
                      data=df,
                      importance=TRUE,
                      ntree=1000)
#Displaying the model
fitNew
```

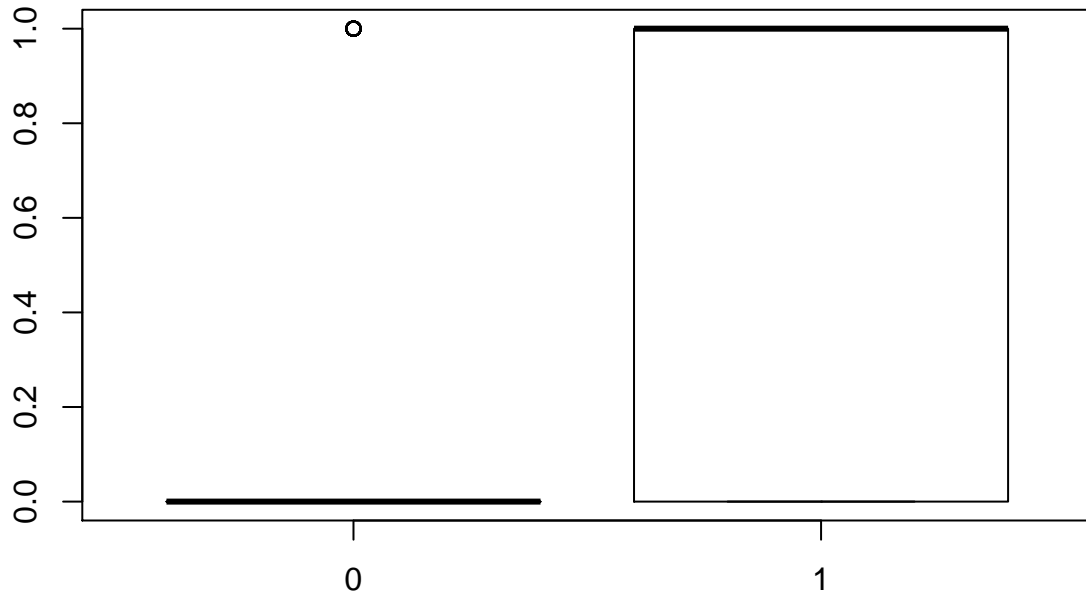
```
##
## Call:
## randomForest(formula = as.factor(survived) ~ pclass + Title + sex + sibsp + fare, data =
##               Type of random forest: classification
##               Number of trees: 1000
## No. of variables tried at each split: 2
##
## OOB estimate of error rate: 18.74%
## Confusion matrix:
##      0   1 class.error
## 0 557  80  0.1255887
## 1 116 293  0.2836186
```

```
#plotting
varImpPlot(fitNew)
```



```
#predicting the test set
yhatRandom1<-predict(fitNew,newdata=testNew)
yhatRandomROC1<-predict(fitNew,newdata=testNew,type="prob")[,2]

#plotting the Random Forest
plot(yhatRandom1,testNew$survived)
```



```
#percentage of correct predictions
mean(yhatRandom1==testNew$survived)
```

```
## [1] 0.8015267
```

- (c) Compare the accuracy of each of the models from this problem set using ROC curves. Comment on which statistical learning method works best for predicting survival of the Titanic passengers.

```
par(mfrow=c(2,2))
g1 <- roc(survived~yhat, data = test)
plot(g1,print.auc=TRUE, main="1st Logistic regression model")
```

```
##
## Call:
## roc.formula(formula = survived ~ yhat, data = test)
##
## Data: yhat in 171 controls (survived 0) < 91 cases (survived 1).
## Area under the curve: 0.6828
```

```
g2 <- roc(survived~yhatlogi, data = testNew)
plot(g2,print.auc=TRUE, main="2nd Logistic regression model")
```

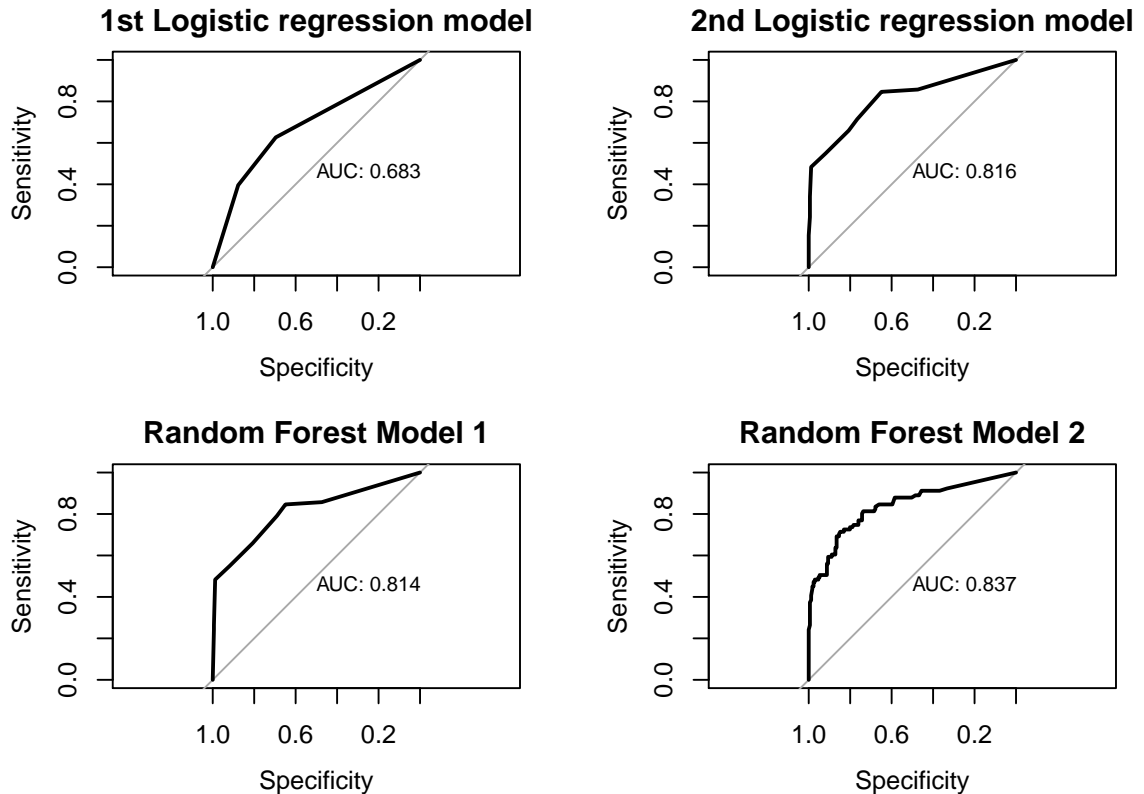
```
##
## Call:
## roc.formula(formula = survived ~ yhatlogi, data = testNew)
##
## Data: yhatlogi in 171 controls (survived 0) < 91 cases (survived 1).
## Area under the curve: 0.8156
```

```
g3 <- roc(survived~yhatRandomROC, data = testNew)
plot(g3,print.auc=TRUE, main="Random Forest Model 1")
```

```
##
## Call:
## roc.formula(formula = survived ~ yhatRandomROC, data = testNew)
##
```

```
## Data: yhatRandomROC in 171 controls (survived 0) < 91 cases (survived 1).
## Area under the curve: 0.8139
```

```
g4 <- roc(survived~yhatRandomROC1, data = testNew)
plot(g4, print.auc=TRUE, main="Random Forest Model 2")
```



```
##
## Call:
## roc.formula(formula = survived ~ yhatRandomROC1, data = testNew)
##
## Data: yhatRandomROC1 in 171 controls (survived 0) < 91 cases (survived 1).
## Area under the curve: 0.8365
```

```
par(mfrow=c(1,1))
```

- Finally, we will explore a gradient boosted tree model, using the `xgboost` package written by your fellow UW student Tianqi Chen. `xgboost` stands for “Extreme Gradient Boosting”, which is state-of-the-art software known for its fast training time and predictive accuracy.

- The XGB algorithm can only handle numeric data, so first we need to convert all categorical variables to a different representation, such as a sparse matrix.

```
library(Matrix)
sparse.matrix.train <- sparse.model.matrix(survived~pclass + sex + Title -1, data = trainNew)
sparse.matrix.test <- sparse.model.matrix(survived~pclass + sex + Title -1, data = testNew)
output_vector = train$survived #output vector to be predicted
```

- The following code fits a boosted tree model and produces a plot. Run the code and provide an explanation of the resulting plot.

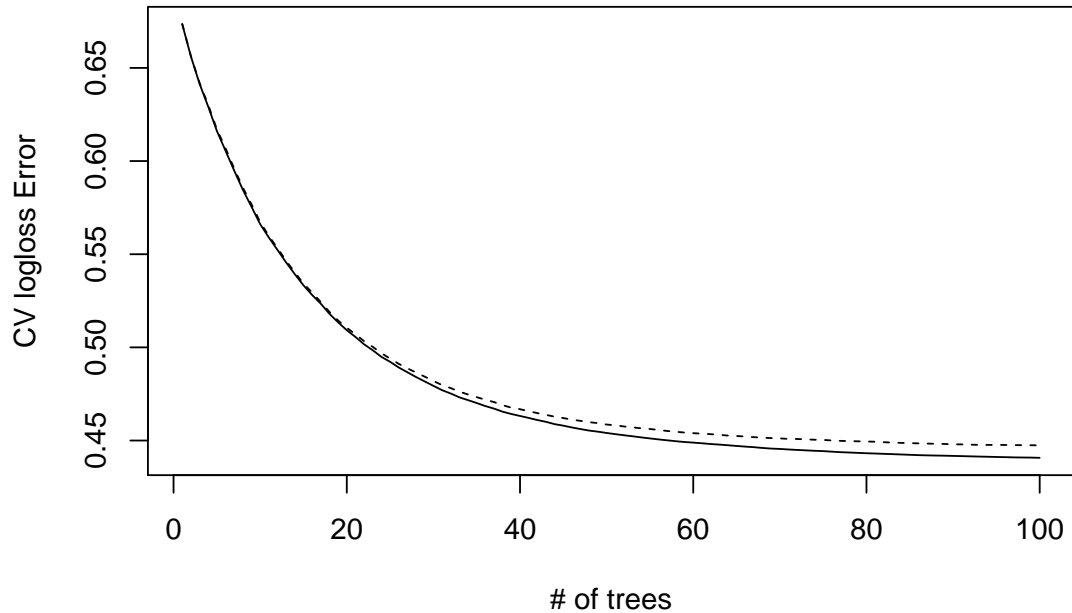

```
xgb.model.one <- xgb.cv(data= sparse.matrix.train,      #train sparse matrix
                        label= output_vector,          #output vector to be predicted
                        eval.metric = 'logloss',        #model minimizes Root Mean Squared Error
                        objective = "reg:logistic",     #regression
                        nfold = 10,
                        #tuning parameters
                        max.depth = 3,                  #Vary btwn 3-15
                        eta = 0.05,                    #Vary btwn 0.1-0.3
                        nthread = 5,                   #Increase this to improve speed
                        subsample= 1,                  #Vary btwn 0.8-1
                        colsample_bytree = 0.5,        #Vary btwn 0.3-0.8
                        lambda = 0.5,                 #Vary between 0-3
                        alpha = 0.5,                  #Vary between 0-3
                        min_child_weight = 3,          #Vary btwn 1-10
                        nround = 100                   #Vary btwn 100-3000 based on max.depth, eta
                        )
```

```
## [0] train-logloss:0.673477+0.000891 test-logloss:0.673599+0.002780
## [1] train-logloss:0.656208+0.002017 test-logloss:0.656367+0.005130
## [2] train-logloss:0.641490+0.003319 test-logloss:0.642021+0.006925
## [3] train-logloss:0.629253+0.006673 test-logloss:0.630212+0.008977
## [4] train-logloss:0.616167+0.006970 test-logloss:0.617088+0.010387
## [5] train-logloss:0.605550+0.007449 test-logloss:0.606916+0.011628
## [6] train-logloss:0.594890+0.006112 test-logloss:0.596016+0.013686
## [7] train-logloss:0.584623+0.006281 test-logloss:0.585790+0.014634
## [8] train-logloss:0.575479+0.005929 test-logloss:0.576785+0.015918
## [9] train-logloss:0.566133+0.006389 test-logloss:0.567253+0.017390
## [10] train-logloss:0.558978+0.006359 test-logloss:0.559819+0.019205
## [11] train-logloss:0.552290+0.006332 test-logloss:0.553257+0.020500
## [12] train-logloss:0.545571+0.006482 test-logloss:0.546488+0.021225
## [13] train-logloss:0.539284+0.006210 test-logloss:0.540277+0.023055
## [14] train-logloss:0.533309+0.005688 test-logloss:0.534235+0.024683
## [15] train-logloss:0.528108+0.006572 test-logloss:0.529408+0.026258
## [16] train-logloss:0.523422+0.005998 test-logloss:0.524337+0.027417
## [17] train-logloss:0.518006+0.006004 test-logloss:0.519025+0.028648
## [18] train-logloss:0.513449+0.005717 test-logloss:0.514431+0.029196
## [19] train-logloss:0.509210+0.005873 test-logloss:0.510510+0.029630
## [20] train-logloss:0.505597+0.005852 test-logloss:0.507015+0.030797
## [21] train-logloss:0.501633+0.006125 test-logloss:0.503316+0.031284
## [22] train-logloss:0.498443+0.006385 test-logloss:0.500248+0.032068
## [23] train-logloss:0.494905+0.006226 test-logloss:0.496850+0.033009
## [24] train-logloss:0.492256+0.006051 test-logloss:0.493980+0.034211
## [25] train-logloss:0.489191+0.006050 test-logloss:0.491156+0.035260
## [26] train-logloss:0.486630+0.006183 test-logloss:0.488738+0.036158
## [27] train-logloss:0.484068+0.006351 test-logloss:0.486505+0.036627
## [28] train-logloss:0.481712+0.006479 test-logloss:0.484283+0.037200
## [29] train-logloss:0.479363+0.006605 test-logloss:0.481953+0.037869
## [30] train-logloss:0.477036+0.006700 test-logloss:0.479764+0.038532
## [31] train-logloss:0.475242+0.006785 test-logloss:0.478132+0.039190
## [32] train-logloss:0.473195+0.006804 test-logloss:0.476246+0.040088
## [33] train-logloss:0.471669+0.006603 test-logloss:0.474709+0.040594
## [34] train-logloss:0.470133+0.006545 test-logloss:0.473275+0.041348
## [35] train-logloss:0.468497+0.006751 test-logloss:0.471879+0.041753
## [36] train-logloss:0.467122+0.006833 test-logloss:0.470630+0.042332
```

```
## [37] train-logloss:0.465471+0.006894 test-logloss:0.468959+0.042888
## [38] train-logloss:0.464232+0.006671 test-logloss:0.467629+0.043468
## [39] train-logloss:0.463182+0.006724 test-logloss:0.466811+0.043834
## [40] train-logloss:0.462126+0.006631 test-logloss:0.465739+0.044498
## [41] train-logloss:0.461008+0.006565 test-logloss:0.464725+0.045192
## [42] train-logloss:0.460032+0.006442 test-logloss:0.463776+0.045883
## [43] train-logloss:0.458790+0.006544 test-logloss:0.462794+0.046327
## [44] train-logloss:0.457982+0.006488 test-logloss:0.462055+0.046768
## [45] train-logloss:0.457009+0.006549 test-logloss:0.461207+0.047194
## [46] train-logloss:0.456140+0.006514 test-logloss:0.460543+0.047749
## [47] train-logloss:0.455331+0.006408 test-logloss:0.459857+0.048213
## [48] train-logloss:0.454749+0.006449 test-logloss:0.459326+0.048435
## [49] train-logloss:0.454059+0.006405 test-logloss:0.458584+0.048869
## [50] train-logloss:0.453424+0.006591 test-logloss:0.458143+0.049070
## [51] train-logloss:0.452834+0.006469 test-logloss:0.457541+0.049491
## [52] train-logloss:0.452275+0.006495 test-logloss:0.457070+0.049769
## [53] train-logloss:0.451683+0.006506 test-logloss:0.456541+0.050139
## [54] train-logloss:0.451124+0.006626 test-logloss:0.456176+0.050414
## [55] train-logloss:0.450606+0.006498 test-logloss:0.455689+0.050728
## [56] train-logloss:0.450102+0.006382 test-logloss:0.455122+0.051161
## [57] train-logloss:0.449633+0.006365 test-logloss:0.454716+0.051495
## [58] train-logloss:0.449203+0.006220 test-logloss:0.454265+0.052011
## [59] train-logloss:0.448882+0.006227 test-logloss:0.453946+0.052179
## [60] train-logloss:0.448522+0.006275 test-logloss:0.453694+0.052482
## [61] train-logloss:0.448130+0.006312 test-logloss:0.453410+0.052592
## [62] train-logloss:0.447812+0.006382 test-logloss:0.453056+0.052859
## [63] train-logloss:0.447461+0.006340 test-logloss:0.452664+0.053086
## [64] train-logloss:0.447088+0.006362 test-logloss:0.452422+0.053292
## [65] train-logloss:0.446740+0.006416 test-logloss:0.452154+0.053392
## [66] train-logloss:0.446425+0.006434 test-logloss:0.451803+0.053585
## [67] train-logloss:0.446012+0.006457 test-logloss:0.451548+0.053864
## [68] train-logloss:0.445680+0.006468 test-logloss:0.451309+0.054063
## [69] train-logloss:0.445481+0.006503 test-logloss:0.451142+0.054146
## [70] train-logloss:0.445193+0.006496 test-logloss:0.450864+0.054364
## [71] train-logloss:0.444945+0.006512 test-logloss:0.450822+0.054549
## [72] train-logloss:0.444721+0.006527 test-logloss:0.450626+0.054634
## [73] train-logloss:0.444497+0.006493 test-logloss:0.450424+0.054808
## [74] train-logloss:0.444310+0.006526 test-logloss:0.450320+0.054874
## [75] train-logloss:0.444019+0.006464 test-logloss:0.450033+0.055150
## [76] train-logloss:0.443751+0.006534 test-logloss:0.449812+0.055298
## [77] train-logloss:0.443564+0.006477 test-logloss:0.449686+0.055431
## [78] train-logloss:0.443366+0.006499 test-logloss:0.449573+0.055535
## [79] train-logloss:0.443167+0.006532 test-logloss:0.449467+0.055743
## [80] train-logloss:0.443030+0.006510 test-logloss:0.449349+0.055888
## [81] train-logloss:0.442861+0.006404 test-logloss:0.449096+0.056079
## [82] train-logloss:0.442680+0.006359 test-logloss:0.448892+0.056151
## [83] train-logloss:0.442522+0.006375 test-logloss:0.448735+0.056158
## [84] train-logloss:0.442337+0.006395 test-logloss:0.448604+0.056293
## [85] train-logloss:0.442158+0.006498 test-logloss:0.448475+0.056306
## [86] train-logloss:0.442028+0.006508 test-logloss:0.448352+0.056440
## [87] train-logloss:0.441942+0.006482 test-logloss:0.448245+0.056506
## [88] train-logloss:0.441856+0.006494 test-logloss:0.448150+0.056645
## [89] train-logloss:0.441736+0.006450 test-logloss:0.447981+0.056810
## [90] train-logloss:0.441640+0.006498 test-logloss:0.447926+0.056842
```

```
## [91] train-logloss:0.441504+0.006451 test-logloss:0.447785+0.056890
## [92] train-logloss:0.441385+0.006479 test-logloss:0.447670+0.056972
## [93] train-logloss:0.441292+0.006499 test-logloss:0.447633+0.057057
## [94] train-logloss:0.441168+0.006562 test-logloss:0.447661+0.057174
## [95] train-logloss:0.441071+0.006556 test-logloss:0.447612+0.057274
## [96] train-logloss:0.440957+0.006583 test-logloss:0.447529+0.057312
## [97] train-logloss:0.440892+0.006591 test-logloss:0.447516+0.057381
## [98] train-logloss:0.440832+0.006602 test-logloss:0.447463+0.057453
## [99] train-logloss:0.440722+0.006587 test-logloss:0.447425+0.057551
```

```
plot(data.frame(xgb.model.one)[,1], type='l', col='black', ylab='CV logloss Error', xlab='# of trees')
lines(data.frame(xgb.model.one)[,3], type='l', lty=2, col='black')
```



```
xgb.model.two <- xgb.cv(data= sparse.matrix.train,      #train sparse matrix
                        label= output_vector,          #output vector to be predicted
                        eval.metric = 'logloss',        #model minimizes Root Mean Squared Error
                        objective = "reg:logistic",     #regression
                        nfold = 10,
                        #tuning parameters
                        max.depth = 3,                  #Vary btwn 3-15
                        eta = 0.05,                     #Vary btwn 0.1-0.3
                        nthread = 5,                    #Increase this to improve speed
                        subsample= 1,                   #Vary btwn 0.8-1
                        colsample_bytree = 0.5,         #Vary btwn 0.3-0.8
                        lambda = 0.5,                  #Vary between 0-3
                        alpha = 0.5,                   #Vary between 0-3
                        min_child_weight = 3,           #Vary btwn 1-10
                        nround = 500                    #Vary btwn 100-3000 based on max.depth, eta
                        )
```

```
## [0] train-logloss:0.676016+0.003033 test-logloss:0.675914+0.004785
## [1] train-logloss:0.660772+0.003124 test-logloss:0.660633+0.006301
## [2] train-logloss:0.646001+0.003957 test-logloss:0.645958+0.008002
## [3] train-logloss:0.631414+0.005456 test-logloss:0.631559+0.009115
## [4] train-logloss:0.618403+0.005416 test-logloss:0.618695+0.009544
## [5] train-logloss:0.605788+0.005540 test-logloss:0.606463+0.011707
```

```
## [6] train-logloss:0.594367+0.005479 test-logloss:0.595402+0.013421
## [7] train-logloss:0.585110+0.004692 test-logloss:0.586636+0.014666
## [8] train-logloss:0.576500+0.005714 test-logloss:0.578187+0.016673
## [9] train-logloss:0.568812+0.005457 test-logloss:0.570519+0.017864
## [10] train-logloss:0.560707+0.005616 test-logloss:0.562416+0.019885
## [11] train-logloss:0.553742+0.006750 test-logloss:0.555787+0.020264
## [12] train-logloss:0.546772+0.005617 test-logloss:0.549075+0.021399
## [13] train-logloss:0.540064+0.005672 test-logloss:0.542731+0.022720
## [14] train-logloss:0.533818+0.005936 test-logloss:0.536846+0.023594
## [15] train-logloss:0.528798+0.006376 test-logloss:0.532112+0.024603
## [16] train-logloss:0.522904+0.006868 test-logloss:0.526416+0.025676
## [17] train-logloss:0.518235+0.006683 test-logloss:0.521843+0.026949
## [18] train-logloss:0.513541+0.007049 test-logloss:0.517137+0.027942
## [19] train-logloss:0.509292+0.007185 test-logloss:0.512908+0.029084
## [20] train-logloss:0.505468+0.006902 test-logloss:0.509006+0.030430
## [21] train-logloss:0.501872+0.006525 test-logloss:0.505334+0.031886
## [22] train-logloss:0.498275+0.006352 test-logloss:0.501790+0.032908
## [23] train-logloss:0.495264+0.005815 test-logloss:0.498621+0.034402
## [24] train-logloss:0.492413+0.005877 test-logloss:0.495535+0.035361
## [25] train-logloss:0.489500+0.006223 test-logloss:0.492665+0.036174
## [26] train-logloss:0.486979+0.006344 test-logloss:0.490203+0.037277
## [27] train-logloss:0.484769+0.006286 test-logloss:0.488088+0.038329
## [28] train-logloss:0.482692+0.006569 test-logloss:0.485800+0.039073
## [29] train-logloss:0.480636+0.006723 test-logloss:0.484112+0.039646
## [30] train-logloss:0.478439+0.006506 test-logloss:0.481989+0.040688
## [31] train-logloss:0.476677+0.006335 test-logloss:0.480352+0.041504
## [32] train-logloss:0.474933+0.006362 test-logloss:0.478489+0.042279
## [33] train-logloss:0.473189+0.006411 test-logloss:0.476978+0.042956
## [34] train-logloss:0.471481+0.006467 test-logloss:0.475465+0.043813
## [35] train-logloss:0.469808+0.006254 test-logloss:0.473865+0.044515
## [36] train-logloss:0.468321+0.006486 test-logloss:0.472646+0.045151
## [37] train-logloss:0.467043+0.006273 test-logloss:0.471353+0.045842
## [38] train-logloss:0.465695+0.006222 test-logloss:0.470191+0.046427
## [39] train-logloss:0.464440+0.006408 test-logloss:0.469110+0.046732
## [40] train-logloss:0.463246+0.006493 test-logloss:0.467972+0.047042
## [41] train-logloss:0.462316+0.006527 test-logloss:0.466976+0.047657
## [42] train-logloss:0.461380+0.006523 test-logloss:0.466034+0.048232
## [43] train-logloss:0.460588+0.006384 test-logloss:0.465311+0.048767
## [44] train-logloss:0.459641+0.006492 test-logloss:0.464379+0.049355
## [45] train-logloss:0.458957+0.006467 test-logloss:0.463716+0.049901
## [46] train-logloss:0.458041+0.006417 test-logloss:0.462900+0.050373
## [47] train-logloss:0.457305+0.006383 test-logloss:0.462052+0.050660
## [48] train-logloss:0.456409+0.006299 test-logloss:0.461193+0.051282
## [49] train-logloss:0.455706+0.006326 test-logloss:0.460517+0.051667
## [50] train-logloss:0.455140+0.006355 test-logloss:0.459948+0.052085
## [51] train-logloss:0.454485+0.006524 test-logloss:0.459409+0.052344
## [52] train-logloss:0.453810+0.006339 test-logloss:0.458780+0.052880
## [53] train-logloss:0.453161+0.006378 test-logloss:0.458348+0.053099
## [54] train-logloss:0.452514+0.006285 test-logloss:0.457724+0.053484
## [55] train-logloss:0.451903+0.006317 test-logloss:0.457191+0.053910
## [56] train-logloss:0.451464+0.006200 test-logloss:0.456725+0.054246
## [57] train-logloss:0.450976+0.006239 test-logloss:0.456268+0.054403
## [58] train-logloss:0.450425+0.006225 test-logloss:0.455835+0.054649
## [59] train-logloss:0.449852+0.006356 test-logloss:0.455415+0.054742
```

```

## [60] train-logloss:0.449360+0.006414 test-logloss:0.455064+0.054894
## [61] train-logloss:0.448934+0.006366 test-logloss:0.454775+0.055099
## [62] train-logloss:0.448459+0.006318 test-logloss:0.454373+0.055381
## [63] train-logloss:0.448207+0.006310 test-logloss:0.454201+0.055499
## [64] train-logloss:0.447708+0.006297 test-logloss:0.453743+0.055847
## [65] train-logloss:0.447391+0.006305 test-logloss:0.453269+0.056070
## [66] train-logloss:0.447089+0.006185 test-logloss:0.452872+0.056386
## [67] train-logloss:0.446827+0.006170 test-logloss:0.452498+0.056631
## [68] train-logloss:0.446470+0.006066 test-logloss:0.452112+0.056926
## [69] train-logloss:0.446195+0.006042 test-logloss:0.451911+0.057085
## [70] train-logloss:0.445950+0.006131 test-logloss:0.451761+0.057179
## [71] train-logloss:0.445669+0.006143 test-logloss:0.451430+0.057307
## [72] train-logloss:0.445470+0.006132 test-logloss:0.451341+0.057472
## [73] train-logloss:0.445227+0.006206 test-logloss:0.451170+0.057524
## [74] train-logloss:0.445084+0.006238 test-logloss:0.451128+0.057567
## [75] train-logloss:0.444840+0.006194 test-logloss:0.450872+0.057697
## [76] train-logloss:0.444575+0.006282 test-logloss:0.450754+0.057775
## [77] train-logloss:0.444415+0.006250 test-logloss:0.450537+0.057928
## [78] train-logloss:0.444234+0.006249 test-logloss:0.450455+0.058001
## [79] train-logloss:0.443978+0.006247 test-logloss:0.450287+0.058158
## [80] train-logloss:0.443813+0.006277 test-logloss:0.450167+0.058255
## [81] train-logloss:0.443607+0.006373 test-logloss:0.450059+0.058296
## [82] train-logloss:0.443345+0.006445 test-logloss:0.449767+0.058374
## [83] train-logloss:0.443177+0.006423 test-logloss:0.449669+0.058424
## [84] train-logloss:0.443007+0.006320 test-logloss:0.449490+0.058713
## [85] train-logloss:0.442895+0.006295 test-logloss:0.449452+0.058790
## [86] train-logloss:0.442769+0.006256 test-logloss:0.449385+0.058986
## [87] train-logloss:0.442581+0.006298 test-logloss:0.449214+0.059087
## [88] train-logloss:0.442407+0.006325 test-logloss:0.449185+0.059027
## [89] train-logloss:0.442215+0.006381 test-logloss:0.449117+0.058995
## [90] train-logloss:0.442090+0.006388 test-logloss:0.448979+0.059042
## [91] train-logloss:0.441953+0.006467 test-logloss:0.448911+0.058991
## [92] train-logloss:0.441867+0.006477 test-logloss:0.448859+0.058967
## [93] train-logloss:0.441797+0.006474 test-logloss:0.448816+0.059046
## [94] train-logloss:0.441685+0.006529 test-logloss:0.448767+0.059073
## [95] train-logloss:0.441568+0.006574 test-logloss:0.448599+0.059016
## [96] train-logloss:0.441475+0.006560 test-logloss:0.448484+0.059063
## [97] train-logloss:0.441395+0.006558 test-logloss:0.448387+0.059103
## [98] train-logloss:0.441298+0.006538 test-logloss:0.448354+0.059210
## [99] train-logloss:0.441201+0.006572 test-logloss:0.448281+0.059281
## [100] train-logloss:0.441100+0.006578 test-logloss:0.448284+0.059353
## [101] train-logloss:0.441049+0.006581 test-logloss:0.448285+0.059410
## [102] train-logloss:0.440939+0.006546 test-logloss:0.448187+0.059582
## [103] train-logloss:0.440837+0.006594 test-logloss:0.448066+0.059532
## [104] train-logloss:0.440715+0.006591 test-logloss:0.447927+0.059586
## [105] train-logloss:0.440654+0.006583 test-logloss:0.447959+0.059656
## [106] train-logloss:0.440605+0.006577 test-logloss:0.447894+0.059723
## [107] train-logloss:0.440547+0.006567 test-logloss:0.447863+0.059809
## [108] train-logloss:0.440476+0.006568 test-logloss:0.447805+0.059833
## [109] train-logloss:0.440412+0.006582 test-logloss:0.447777+0.059824
## [110] train-logloss:0.440339+0.006558 test-logloss:0.447693+0.059910
## [111] train-logloss:0.440283+0.006518 test-logloss:0.447694+0.059900
## [112] train-logloss:0.440205+0.006492 test-logloss:0.447629+0.059949
## [113] train-logloss:0.440128+0.006508 test-logloss:0.447597+0.060010

```

```

## [114] train-logloss:0.440064+0.006488 test-logloss:0.447572+0.060064
## [115] train-logloss:0.439975+0.006428 test-logloss:0.447547+0.060115
## [116] train-logloss:0.439885+0.006425 test-logloss:0.447436+0.060170
## [117] train-logloss:0.439827+0.006451 test-logloss:0.447429+0.060174
## [118] train-logloss:0.439788+0.006453 test-logloss:0.447411+0.060176
## [119] train-logloss:0.439734+0.006477 test-logloss:0.447431+0.060175
## [120] train-logloss:0.439654+0.006493 test-logloss:0.447375+0.060143
## [121] train-logloss:0.439633+0.006486 test-logloss:0.447348+0.060184
## [122] train-logloss:0.439591+0.006478 test-logloss:0.447283+0.060249
## [123] train-logloss:0.439566+0.006483 test-logloss:0.447314+0.060255
## [124] train-logloss:0.439502+0.006488 test-logloss:0.447315+0.060276
## [125] train-logloss:0.439466+0.006507 test-logloss:0.447299+0.060210
## [126] train-logloss:0.439404+0.006506 test-logloss:0.447226+0.060211
## [127] train-logloss:0.439319+0.006472 test-logloss:0.447146+0.060243
## [128] train-logloss:0.439296+0.006473 test-logloss:0.447130+0.060257
## [129] train-logloss:0.439272+0.006463 test-logloss:0.447135+0.060283
## [130] train-logloss:0.439227+0.006467 test-logloss:0.447079+0.060294
## [131] train-logloss:0.439173+0.006473 test-logloss:0.447043+0.060327
## [132] train-logloss:0.439132+0.006467 test-logloss:0.447035+0.060376
## [133] train-logloss:0.439085+0.006471 test-logloss:0.446997+0.060413
## [134] train-logloss:0.439059+0.006477 test-logloss:0.446963+0.060403
## [135] train-logloss:0.439018+0.006481 test-logloss:0.446928+0.060411
## [136] train-logloss:0.438984+0.006488 test-logloss:0.446895+0.060392
## [137] train-logloss:0.438945+0.006462 test-logloss:0.446861+0.060489
## [138] train-logloss:0.438899+0.006485 test-logloss:0.446817+0.060440
## [139] train-logloss:0.438871+0.006483 test-logloss:0.446793+0.060476
## [140] train-logloss:0.438833+0.006490 test-logloss:0.446740+0.060482
## [141] train-logloss:0.438790+0.006511 test-logloss:0.446719+0.060479
## [142] train-logloss:0.438745+0.006567 test-logloss:0.446700+0.060438
## [143] train-logloss:0.438689+0.006560 test-logloss:0.446656+0.060499
## [144] train-logloss:0.438668+0.006554 test-logloss:0.446640+0.060532
## [145] train-logloss:0.438614+0.006582 test-logloss:0.446620+0.060514
## [146] train-logloss:0.438562+0.006579 test-logloss:0.446590+0.060554
## [147] train-logloss:0.438530+0.006569 test-logloss:0.446593+0.060576
## [148] train-logloss:0.438520+0.006572 test-logloss:0.446604+0.060567
## [149] train-logloss:0.438470+0.006585 test-logloss:0.446583+0.060589
## [150] train-logloss:0.438408+0.006621 test-logloss:0.446514+0.060577
## [151] train-logloss:0.438379+0.006626 test-logloss:0.446485+0.060585
## [152] train-logloss:0.438345+0.006616 test-logloss:0.446472+0.060600
## [153] train-logloss:0.438299+0.006630 test-logloss:0.446436+0.060575
## [154] train-logloss:0.438247+0.006636 test-logloss:0.446431+0.060579
## [155] train-logloss:0.438196+0.006650 test-logloss:0.446445+0.060499
## [156] train-logloss:0.438156+0.006635 test-logloss:0.446457+0.060467
## [157] train-logloss:0.438139+0.006632 test-logloss:0.446439+0.060480
## [158] train-logloss:0.438101+0.006657 test-logloss:0.446407+0.060455
## [159] train-logloss:0.438084+0.006664 test-logloss:0.446401+0.060471
## [160] train-logloss:0.438048+0.006644 test-logloss:0.446395+0.060455
## [161] train-logloss:0.438020+0.006656 test-logloss:0.446361+0.060428
## [162] train-logloss:0.437994+0.006642 test-logloss:0.446354+0.060467
## [163] train-logloss:0.437972+0.006641 test-logloss:0.446330+0.060469
## [164] train-logloss:0.437962+0.006642 test-logloss:0.446331+0.060457
## [165] train-logloss:0.437940+0.006642 test-logloss:0.446298+0.060458
## [166] train-logloss:0.437920+0.006636 test-logloss:0.446244+0.060473
## [167] train-logloss:0.437901+0.006628 test-logloss:0.446235+0.060450

```

```

## [168] train-logloss:0.437879+0.006613 test-logloss:0.446227+0.060498
## [169] train-logloss:0.437861+0.006610 test-logloss:0.446217+0.060450
## [170] train-logloss:0.437834+0.006617 test-logloss:0.446252+0.060420
## [171] train-logloss:0.437824+0.006620 test-logloss:0.446278+0.060425
## [172] train-logloss:0.437797+0.006645 test-logloss:0.446300+0.060425
## [173] train-logloss:0.437773+0.006639 test-logloss:0.446235+0.060439
## [174] train-logloss:0.437722+0.006652 test-logloss:0.446213+0.060410
## [175] train-logloss:0.437698+0.006657 test-logloss:0.446199+0.060435
## [176] train-logloss:0.437677+0.006662 test-logloss:0.446167+0.060428
## [177] train-logloss:0.437654+0.006656 test-logloss:0.446119+0.060459
## [178] train-logloss:0.437630+0.006651 test-logloss:0.446101+0.060506
## [179] train-logloss:0.437615+0.006651 test-logloss:0.446101+0.060522
## [180] train-logloss:0.437592+0.006666 test-logloss:0.446084+0.060517
## [181] train-logloss:0.437558+0.006683 test-logloss:0.446029+0.060468
## [182] train-logloss:0.437549+0.006684 test-logloss:0.446026+0.060463
## [183] train-logloss:0.437541+0.006682 test-logloss:0.446018+0.060455
## [184] train-logloss:0.437518+0.006679 test-logloss:0.446038+0.060415
## [185] train-logloss:0.437511+0.006678 test-logloss:0.446024+0.060410
## [186] train-logloss:0.437503+0.006673 test-logloss:0.446007+0.060411
## [187] train-logloss:0.437497+0.006673 test-logloss:0.446004+0.060414
## [188] train-logloss:0.437476+0.006663 test-logloss:0.445982+0.060448
## [189] train-logloss:0.437459+0.006662 test-logloss:0.445986+0.060442
## [190] train-logloss:0.437438+0.006650 test-logloss:0.445973+0.060459
## [191] train-logloss:0.437415+0.006636 test-logloss:0.445924+0.060496
## [192] train-logloss:0.437389+0.006651 test-logloss:0.445900+0.060486
## [193] train-logloss:0.437363+0.006650 test-logloss:0.445852+0.060521
## [194] train-logloss:0.437353+0.006653 test-logloss:0.445847+0.060519
## [195] train-logloss:0.437344+0.006651 test-logloss:0.445827+0.060525
## [196] train-logloss:0.437330+0.006640 test-logloss:0.445797+0.060598
## [197] train-logloss:0.437304+0.006653 test-logloss:0.445783+0.060604
## [198] train-logloss:0.437295+0.006651 test-logloss:0.445784+0.060595
## [199] train-logloss:0.437270+0.006657 test-logloss:0.445753+0.060591
## [200] train-logloss:0.437258+0.006655 test-logloss:0.445749+0.060585
## [201] train-logloss:0.437233+0.006660 test-logloss:0.445752+0.060575
## [202] train-logloss:0.437226+0.006658 test-logloss:0.445734+0.060571
## [203] train-logloss:0.437217+0.006654 test-logloss:0.445731+0.060594
## [204] train-logloss:0.437192+0.006636 test-logloss:0.445703+0.060586
## [205] train-logloss:0.437184+0.006639 test-logloss:0.445683+0.060569
## [206] train-logloss:0.437161+0.006642 test-logloss:0.445661+0.060582
## [207] train-logloss:0.437156+0.006643 test-logloss:0.445657+0.060570
## [208] train-logloss:0.437138+0.006637 test-logloss:0.445648+0.060591
## [209] train-logloss:0.437134+0.006635 test-logloss:0.445646+0.060585
## [210] train-logloss:0.437120+0.006633 test-logloss:0.445626+0.060576
## [211] train-logloss:0.437101+0.006632 test-logloss:0.445592+0.060577
## [212] train-logloss:0.437080+0.006618 test-logloss:0.445596+0.060574
## [213] train-logloss:0.437072+0.006619 test-logloss:0.445592+0.060568
## [214] train-logloss:0.437066+0.006618 test-logloss:0.445598+0.060548
## [215] train-logloss:0.437057+0.006617 test-logloss:0.445585+0.060552
## [216] train-logloss:0.437045+0.006616 test-logloss:0.445586+0.060554
## [217] train-logloss:0.437027+0.006621 test-logloss:0.445582+0.060561
## [218] train-logloss:0.437019+0.006620 test-logloss:0.445576+0.060557
## [219] train-logloss:0.437015+0.006621 test-logloss:0.445585+0.060541
## [220] train-logloss:0.437008+0.006622 test-logloss:0.445595+0.060526
## [221] train-logloss:0.437003+0.006622 test-logloss:0.445582+0.060517

```

```

## [222] train-logloss:0.436984+0.006606 test-logloss:0.445555+0.060555
## [223] train-logloss:0.436972+0.006602 test-logloss:0.445554+0.060600
## [224] train-logloss:0.436965+0.006600 test-logloss:0.445552+0.060588
## [225] train-logloss:0.436959+0.006604 test-logloss:0.445548+0.060593
## [226] train-logloss:0.436936+0.006588 test-logloss:0.445549+0.060604
## [227] train-logloss:0.436927+0.006594 test-logloss:0.445530+0.060605
## [228] train-logloss:0.436901+0.006617 test-logloss:0.445505+0.060593
## [229] train-logloss:0.436885+0.006620 test-logloss:0.445487+0.060604
## [230] train-logloss:0.436880+0.006620 test-logloss:0.445480+0.060592
## [231] train-logloss:0.436874+0.006619 test-logloss:0.445473+0.060585
## [232] train-logloss:0.436868+0.006617 test-logloss:0.445468+0.060580
## [233] train-logloss:0.436862+0.006618 test-logloss:0.445468+0.060580
## [234] train-logloss:0.436857+0.006619 test-logloss:0.445456+0.060570
## [235] train-logloss:0.436845+0.006610 test-logloss:0.445418+0.060583
## [236] train-logloss:0.436837+0.006611 test-logloss:0.445406+0.060576
## [237] train-logloss:0.436822+0.006602 test-logloss:0.445381+0.060572
## [238] train-logloss:0.436809+0.006611 test-logloss:0.445383+0.060543
## [239] train-logloss:0.436806+0.006614 test-logloss:0.445385+0.060534
## [240] train-logloss:0.436802+0.006612 test-logloss:0.445378+0.060536
## [241] train-logloss:0.436796+0.006617 test-logloss:0.445379+0.060530
## [242] train-logloss:0.436788+0.006621 test-logloss:0.445386+0.060527
## [243] train-logloss:0.436785+0.006618 test-logloss:0.445379+0.060531
## [244] train-logloss:0.436775+0.006611 test-logloss:0.445373+0.060542
## [245] train-logloss:0.436764+0.006603 test-logloss:0.445343+0.060549
## [246] train-logloss:0.436754+0.006608 test-logloss:0.445324+0.060535
## [247] train-logloss:0.436743+0.006607 test-logloss:0.445298+0.060516
## [248] train-logloss:0.436726+0.006599 test-logloss:0.445306+0.060502
## [249] train-logloss:0.436722+0.006600 test-logloss:0.445301+0.060500
## [250] train-logloss:0.436718+0.006600 test-logloss:0.445310+0.060492
## [251] train-logloss:0.436710+0.006601 test-logloss:0.445314+0.060474
## [252] train-logloss:0.436702+0.006597 test-logloss:0.445304+0.060468
## [253] train-logloss:0.436693+0.006601 test-logloss:0.445311+0.060454
## [254] train-logloss:0.436686+0.006601 test-logloss:0.445291+0.060440
## [255] train-logloss:0.436681+0.006604 test-logloss:0.445292+0.060429
## [256] train-logloss:0.436667+0.006612 test-logloss:0.445282+0.060432
## [257] train-logloss:0.436652+0.006621 test-logloss:0.445260+0.060423
## [258] train-logloss:0.436641+0.006630 test-logloss:0.445237+0.060424
## [259] train-logloss:0.436631+0.006637 test-logloss:0.445232+0.060413
## [260] train-logloss:0.436628+0.006637 test-logloss:0.445238+0.060404
## [261] train-logloss:0.436615+0.006639 test-logloss:0.445233+0.060402
## [262] train-logloss:0.436609+0.006639 test-logloss:0.445225+0.060397
## [263] train-logloss:0.436607+0.006639 test-logloss:0.445224+0.060389
## [264] train-logloss:0.436601+0.006640 test-logloss:0.445236+0.060382
## [265] train-logloss:0.436593+0.006645 test-logloss:0.445244+0.060385
## [266] train-logloss:0.436584+0.006644 test-logloss:0.445241+0.060366
## [267] train-logloss:0.436582+0.006645 test-logloss:0.445234+0.060365
## [268] train-logloss:0.436577+0.006643 test-logloss:0.445229+0.060364
## [269] train-logloss:0.436565+0.006641 test-logloss:0.445195+0.060373
## [270] train-logloss:0.436561+0.006640 test-logloss:0.445193+0.060372
## [271] train-logloss:0.436559+0.006640 test-logloss:0.445191+0.060372
## [272] train-logloss:0.436558+0.006639 test-logloss:0.445192+0.060371
## [273] train-logloss:0.436550+0.006639 test-logloss:0.445171+0.060370
## [274] train-logloss:0.436538+0.006636 test-logloss:0.445184+0.060375
## [275] train-logloss:0.436530+0.006633 test-logloss:0.445194+0.060377

```



```

## [276] train-logloss:0.436525+0.006636 test-logloss:0.445201+0.060380
## [277] train-logloss:0.436521+0.006634 test-logloss:0.445199+0.060366
## [278] train-logloss:0.436517+0.006633 test-logloss:0.445204+0.060359
## [279] train-logloss:0.436514+0.006633 test-logloss:0.445205+0.060366
## [280] train-logloss:0.436513+0.006632 test-logloss:0.445209+0.060362
## [281] train-logloss:0.436508+0.006633 test-logloss:0.445187+0.060356
## [282] train-logloss:0.436499+0.006625 test-logloss:0.445177+0.060387
## [283] train-logloss:0.436497+0.006625 test-logloss:0.445180+0.060379
## [284] train-logloss:0.436485+0.006626 test-logloss:0.445152+0.060372
## [285] train-logloss:0.436483+0.006623 test-logloss:0.445140+0.060384
## [286] train-logloss:0.436476+0.006627 test-logloss:0.445134+0.060362
## [287] train-logloss:0.436473+0.006627 test-logloss:0.445129+0.060356
## [288] train-logloss:0.436471+0.006625 test-logloss:0.445131+0.060354
## [289] train-logloss:0.436471+0.006625 test-logloss:0.445134+0.060355
## [290] train-logloss:0.436467+0.006623 test-logloss:0.445134+0.060363
## [291] train-logloss:0.436464+0.006620 test-logloss:0.445132+0.060366
## [292] train-logloss:0.436464+0.006620 test-logloss:0.445131+0.060365
## [293] train-logloss:0.436459+0.006618 test-logloss:0.445124+0.060366
## [294] train-logloss:0.436456+0.006617 test-logloss:0.445119+0.060375
## [295] train-logloss:0.436455+0.006616 test-logloss:0.445115+0.060380
## [296] train-logloss:0.436453+0.006615 test-logloss:0.445112+0.060380
## [297] train-logloss:0.436451+0.006613 test-logloss:0.445109+0.060383
## [298] train-logloss:0.436450+0.006613 test-logloss:0.445099+0.060386
## [299] train-logloss:0.436443+0.006604 test-logloss:0.445107+0.060383
## [300] train-logloss:0.436443+0.006604 test-logloss:0.445109+0.060382
## [301] train-logloss:0.436437+0.006596 test-logloss:0.445115+0.060378
## [302] train-logloss:0.436435+0.006596 test-logloss:0.445123+0.060377
## [303] train-logloss:0.436435+0.006596 test-logloss:0.445122+0.060377
## [304] train-logloss:0.436428+0.006596 test-logloss:0.445120+0.060380
## [305] train-logloss:0.436426+0.006596 test-logloss:0.445122+0.060381
## [306] train-logloss:0.436426+0.006596 test-logloss:0.445126+0.060378
## [307] train-logloss:0.436422+0.006592 test-logloss:0.445133+0.060374
## [308] train-logloss:0.436421+0.006591 test-logloss:0.445130+0.060376
## [309] train-logloss:0.436421+0.006591 test-logloss:0.445132+0.060376
## [310] train-logloss:0.436419+0.006590 test-logloss:0.445133+0.060378
## [311] train-logloss:0.436417+0.006590 test-logloss:0.445132+0.060377
## [312] train-logloss:0.436416+0.006590 test-logloss:0.445136+0.060377
## [313] train-logloss:0.436416+0.006590 test-logloss:0.445137+0.060377
## [314] train-logloss:0.436416+0.006590 test-logloss:0.445138+0.060376
## [315] train-logloss:0.436416+0.006590 test-logloss:0.445140+0.060376
## [316] train-logloss:0.436416+0.006590 test-logloss:0.445140+0.060376
## [317] train-logloss:0.436416+0.006589 test-logloss:0.445141+0.060376
## [318] train-logloss:0.436416+0.006589 test-logloss:0.445141+0.060376
## [319] train-logloss:0.436415+0.006589 test-logloss:0.445139+0.060376
## [320] train-logloss:0.436415+0.006589 test-logloss:0.445139+0.060376
## [321] train-logloss:0.436415+0.006589 test-logloss:0.445139+0.060376
## [322] train-logloss:0.436415+0.006589 test-logloss:0.445140+0.060376
## [323] train-logloss:0.436415+0.006589 test-logloss:0.445140+0.060376
## [324] train-logloss:0.436414+0.006589 test-logloss:0.445138+0.060376
## [325] train-logloss:0.436414+0.006589 test-logloss:0.445139+0.060376
## [326] train-logloss:0.436414+0.006589 test-logloss:0.445141+0.060375
## [327] train-logloss:0.436413+0.006589 test-logloss:0.445139+0.060376
## [328] train-logloss:0.436413+0.006589 test-logloss:0.445140+0.060376
## [329] train-logloss:0.436413+0.006589 test-logloss:0.445140+0.060376

```

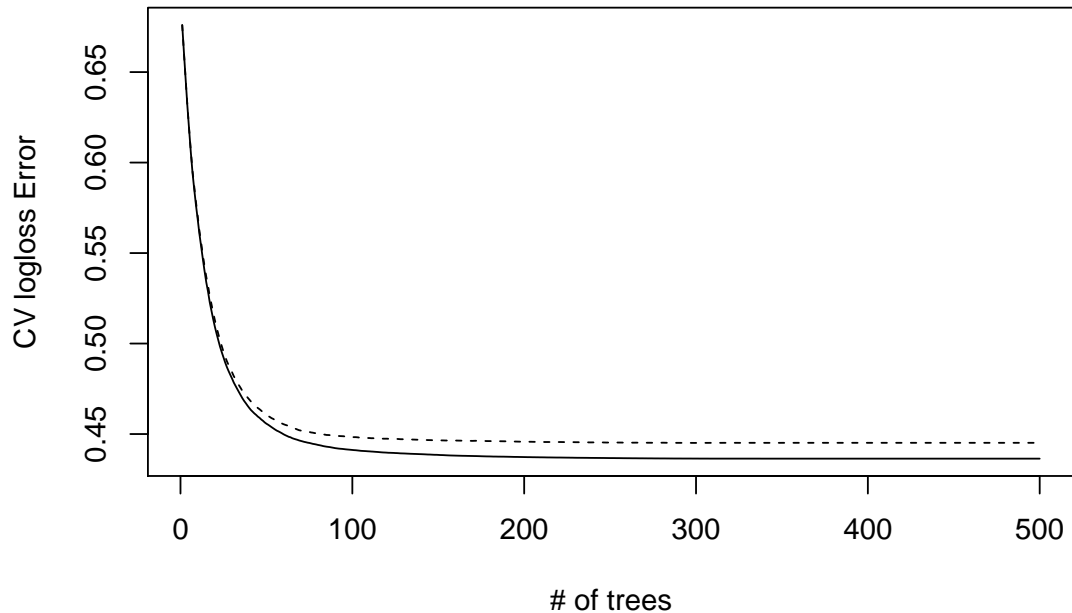
[illegible]

[illegible]

[illegible]

```
## [492] train-logloss:0.436409+0.006588 test-logloss:0.445141+0.060375
## [493] train-logloss:0.436409+0.006588 test-logloss:0.445141+0.060375
## [494] train-logloss:0.436409+0.006588 test-logloss:0.445141+0.060375
## [495] train-logloss:0.436409+0.006588 test-logloss:0.445141+0.060375
## [496] train-logloss:0.436409+0.006588 test-logloss:0.445141+0.060375
## [497] train-logloss:0.436409+0.006588 test-logloss:0.445141+0.060375
## [498] train-logloss:0.436409+0.006588 test-logloss:0.445141+0.060375
## [499] train-logloss:0.436409+0.006588 test-logloss:0.445141+0.060375
```

```
plot(data.frame(xgb.model.two)[,1], type='l', col='black', ylab='CV logloss Error', xlab='# of trees')
lines(data.frame(xgb.model.two)[,3], type='l', lty=2, col='black')
```



The plot produced by the above code displays the behaviour of number of trees vs CV logless error. With the increase in the number of trees, there is decrease in the CV logless error. Initially, there is sudden decrease in the CV logless error with increase in the number of trees but as the number of trees reaches 100, the error remains almost constant. Various paramters in the model are changing with change in the value. Increase in the step size of eta which is the step size shrinkage results into shifting the minimum CV logless error rate towards smaller # of trees. Increase in number of threads increases the speed of the model.

- (c) Modify the code to fit a boosted tree model that allows for 8 levels in each tree and uses a learning rate $\eta = .1$. Produce a visualization comparing the two models and explain what you can conclude about the new model. Which model do you prefer and why?

```
xgb.model.three <- xgb.cv(data= sparse.matrix.train,      #train sparse matrix
                           label= output_vector,        #output vector to be predicted
                           eval.metric = 'logloss',     #model minimizes Root Mean Squared Error
                           objective = "reg:logistic",  #regression
                           nfold = 10,
                           #tuning parameters
                           max.depth = 8,               #Vary btwn 3-15
                           eta = 0.1,                  #Vary btwn 0.1-0.3
                           nthread = 5,                 #Increase this to improve speed
                           subsample= 1,                #Vary btwn 0.8-1
                           colsample_bytree = 0.5,      #Vary btwn 0.3-0.8
                           lambda = 0.5,               #Vary between 0-3
                           alpha = 0.5,                #Vary between 0-3
```

```

min_child_weight = 3,      #Vary btwn 1-10
nround = 500               #Vary btwn 100-3000 based on max.depth, etc
)

```

```

## [0] train-logloss:0.658224+0.007228 test-logloss:0.659521+0.006478
## [1] train-logloss:0.630725+0.007342 test-logloss:0.632418+0.006093
## [2] train-logloss:0.605757+0.008440 test-logloss:0.608125+0.009016
## [3] train-logloss:0.585217+0.006450 test-logloss:0.587661+0.007681
## [4] train-logloss:0.566317+0.008071 test-logloss:0.569675+0.010142
## [5] train-logloss:0.549776+0.006828 test-logloss:0.553678+0.010557
## [6] train-logloss:0.536483+0.006623 test-logloss:0.540846+0.012868
## [7] train-logloss:0.524346+0.006163 test-logloss:0.528827+0.014435
## [8] train-logloss:0.514025+0.005894 test-logloss:0.518265+0.014907
## [9] train-logloss:0.506115+0.006631 test-logloss:0.510776+0.015632
## [10] train-logloss:0.498804+0.007566 test-logloss:0.503598+0.015777
## [11] train-logloss:0.492034+0.007321 test-logloss:0.496605+0.016264
## [12] train-logloss:0.486292+0.007202 test-logloss:0.491569+0.016674
## [13] train-logloss:0.480844+0.007216 test-logloss:0.486386+0.017463
## [14] train-logloss:0.476612+0.006725 test-logloss:0.482195+0.018641
## [15] train-logloss:0.472697+0.006410 test-logloss:0.478711+0.019322
## [16] train-logloss:0.469534+0.006526 test-logloss:0.475304+0.019609
## [17] train-logloss:0.466797+0.006561 test-logloss:0.472331+0.020441
## [18] train-logloss:0.464500+0.005884 test-logloss:0.470375+0.020852
## [19] train-logloss:0.461744+0.005561 test-logloss:0.467816+0.021993
## [20] train-logloss:0.459702+0.005288 test-logloss:0.465754+0.022798
## [21] train-logloss:0.458120+0.004947 test-logloss:0.464001+0.023285
## [22] train-logloss:0.456556+0.004668 test-logloss:0.462084+0.024121
## [23] train-logloss:0.455012+0.004482 test-logloss:0.460550+0.024828
## [24] train-logloss:0.453337+0.004444 test-logloss:0.459006+0.025318
## [25] train-logloss:0.451979+0.004287 test-logloss:0.457587+0.025834
## [26] train-logloss:0.450990+0.004236 test-logloss:0.456682+0.026170
## [27] train-logloss:0.449925+0.004199 test-logloss:0.455686+0.026308
## [28] train-logloss:0.448831+0.004208 test-logloss:0.454738+0.026852
## [29] train-logloss:0.447974+0.004141 test-logloss:0.454076+0.027271
## [30] train-logloss:0.447330+0.004027 test-logloss:0.453593+0.027649
## [31] train-logloss:0.446391+0.003861 test-logloss:0.452623+0.027961
## [32] train-logloss:0.445957+0.003762 test-logloss:0.452317+0.028178
## [33] train-logloss:0.445286+0.003597 test-logloss:0.451788+0.028490
## [34] train-logloss:0.444705+0.003531 test-logloss:0.451389+0.028773
## [35] train-logloss:0.444228+0.003559 test-logloss:0.450819+0.028778
## [36] train-logloss:0.443880+0.003628 test-logloss:0.450585+0.028877
## [37] train-logloss:0.443509+0.003585 test-logloss:0.450120+0.028956
## [38] train-logloss:0.443086+0.003515 test-logloss:0.449738+0.029158
## [39] train-logloss:0.442800+0.003485 test-logloss:0.449605+0.029357
## [40] train-logloss:0.442559+0.003406 test-logloss:0.449351+0.029470
## [41] train-logloss:0.442148+0.003537 test-logloss:0.449107+0.029415
## [42] train-logloss:0.441869+0.003625 test-logloss:0.448951+0.029536
## [43] train-logloss:0.441609+0.003666 test-logloss:0.448718+0.029668
## [44] train-logloss:0.441409+0.003640 test-logloss:0.448564+0.029780
## [45] train-logloss:0.441144+0.003649 test-logloss:0.448289+0.029909
## [46] train-logloss:0.440850+0.003661 test-logloss:0.447860+0.030076
## [47] train-logloss:0.440647+0.003612 test-logloss:0.447695+0.030317
## [48] train-logloss:0.440469+0.003566 test-logloss:0.447611+0.030400
## [49] train-logloss:0.440315+0.003474 test-logloss:0.447372+0.030489

```

```

## [50] train-logloss:0.440121+0.003469 test-logloss:0.447272+0.030522
## [51] train-logloss:0.440004+0.003465 test-logloss:0.447103+0.030509
## [52] train-logloss:0.439956+0.003474 test-logloss:0.447111+0.030618
## [53] train-logloss:0.439816+0.003474 test-logloss:0.447128+0.030641
## [54] train-logloss:0.439695+0.003480 test-logloss:0.447134+0.030773
## [55] train-logloss:0.439550+0.003525 test-logloss:0.447061+0.030810
## [56] train-logloss:0.439450+0.003538 test-logloss:0.446984+0.030834
## [57] train-logloss:0.439325+0.003546 test-logloss:0.446993+0.030981
## [58] train-logloss:0.439206+0.003566 test-logloss:0.447041+0.031065
## [59] train-logloss:0.439078+0.003521 test-logloss:0.447100+0.031067
## [60] train-logloss:0.438931+0.003470 test-logloss:0.446959+0.031082
## [61] train-logloss:0.438836+0.003464 test-logloss:0.446998+0.031124
## [62] train-logloss:0.438715+0.003477 test-logloss:0.446908+0.031074
## [63] train-logloss:0.438654+0.003511 test-logloss:0.446851+0.031110
## [64] train-logloss:0.438573+0.003568 test-logloss:0.446811+0.031069
## [65] train-logloss:0.438479+0.003552 test-logloss:0.446760+0.031094
## [66] train-logloss:0.438426+0.003545 test-logloss:0.446713+0.031204
## [67] train-logloss:0.438358+0.003525 test-logloss:0.446591+0.031308
## [68] train-logloss:0.438325+0.003529 test-logloss:0.446520+0.031306
## [69] train-logloss:0.438273+0.003515 test-logloss:0.446482+0.031354
## [70] train-logloss:0.438173+0.003475 test-logloss:0.446517+0.031377
## [71] train-logloss:0.438152+0.003471 test-logloss:0.446537+0.031378
## [72] train-logloss:0.438082+0.003438 test-logloss:0.446386+0.031498
## [73] train-logloss:0.438006+0.003454 test-logloss:0.446300+0.031459
## [74] train-logloss:0.437978+0.003464 test-logloss:0.446300+0.031486
## [75] train-logloss:0.437954+0.003455 test-logloss:0.446227+0.031556
## [76] train-logloss:0.437913+0.003432 test-logloss:0.446241+0.031569
## [77] train-logloss:0.437866+0.003389 test-logloss:0.446177+0.031616
## [78] train-logloss:0.437831+0.003382 test-logloss:0.446168+0.031622
## [79] train-logloss:0.437776+0.003381 test-logloss:0.446136+0.031702
## [80] train-logloss:0.437729+0.003361 test-logloss:0.446060+0.031731
## [81] train-logloss:0.437684+0.003394 test-logloss:0.446026+0.031731
## [82] train-logloss:0.437630+0.003373 test-logloss:0.446007+0.031760
## [83] train-logloss:0.437569+0.003411 test-logloss:0.445945+0.031790
## [84] train-logloss:0.437542+0.003416 test-logloss:0.445906+0.031819
## [85] train-logloss:0.437509+0.003422 test-logloss:0.445899+0.031820
## [86] train-logloss:0.437445+0.003406 test-logloss:0.445869+0.031884
## [87] train-logloss:0.437393+0.003420 test-logloss:0.445843+0.031840
## [88] train-logloss:0.437381+0.003413 test-logloss:0.445820+0.031820
## [89] train-logloss:0.437341+0.003395 test-logloss:0.445759+0.031912
## [90] train-logloss:0.437298+0.003408 test-logloss:0.445716+0.031893
## [91] train-logloss:0.437268+0.003412 test-logloss:0.445711+0.031887
## [92] train-logloss:0.437246+0.003429 test-logloss:0.445718+0.031917
## [93] train-logloss:0.437232+0.003426 test-logloss:0.445695+0.031937
## [94] train-logloss:0.437209+0.003435 test-logloss:0.445734+0.031942
## [95] train-logloss:0.437183+0.003419 test-logloss:0.445709+0.031944
## [96] train-logloss:0.437153+0.003399 test-logloss:0.445682+0.031968
## [97] train-logloss:0.437130+0.003403 test-logloss:0.445654+0.031964
## [98] train-logloss:0.437117+0.003397 test-logloss:0.445621+0.031963
## [99] train-logloss:0.437113+0.003398 test-logloss:0.445620+0.031955
## [100] train-logloss:0.437081+0.003385 test-logloss:0.445583+0.032010
## [101] train-logloss:0.437067+0.003389 test-logloss:0.445631+0.032058
## [102] train-logloss:0.437028+0.003380 test-logloss:0.445604+0.032090
## [103] train-logloss:0.437013+0.003376 test-logloss:0.445599+0.032096

```

```

## [104] train-logloss:0.436999+0.003381 test-logloss:0.445585+0.032128
## [105] train-logloss:0.436988+0.003387 test-logloss:0.445627+0.032173
## [106] train-logloss:0.436956+0.003401 test-logloss:0.445601+0.032190
## [107] train-logloss:0.436937+0.003393 test-logloss:0.445539+0.032228
## [108] train-logloss:0.436898+0.003359 test-logloss:0.445525+0.032240
## [109] train-logloss:0.436888+0.003367 test-logloss:0.445529+0.032255
## [110] train-logloss:0.436875+0.003358 test-logloss:0.445525+0.032251
## [111] train-logloss:0.436865+0.003358 test-logloss:0.445535+0.032257
## [112] train-logloss:0.436854+0.003360 test-logloss:0.445516+0.032241
## [113] train-logloss:0.436842+0.003358 test-logloss:0.445539+0.032236
## [114] train-logloss:0.436831+0.003362 test-logloss:0.445528+0.032175
## [115] train-logloss:0.436796+0.003366 test-logloss:0.445535+0.032155
## [116] train-logloss:0.436766+0.003362 test-logloss:0.445530+0.032136
## [117] train-logloss:0.436753+0.003367 test-logloss:0.445538+0.032150
## [118] train-logloss:0.436734+0.003361 test-logloss:0.445533+0.032155
## [119] train-logloss:0.436725+0.003362 test-logloss:0.445528+0.032149
## [120] train-logloss:0.436712+0.003373 test-logloss:0.445532+0.032107
## [121] train-logloss:0.436700+0.003379 test-logloss:0.445547+0.032113
## [122] train-logloss:0.436687+0.003387 test-logloss:0.445534+0.032087
## [123] train-logloss:0.436678+0.003395 test-logloss:0.445511+0.032073
## [124] train-logloss:0.436653+0.003400 test-logloss:0.445516+0.032127
## [125] train-logloss:0.436645+0.003409 test-logloss:0.445531+0.032140
## [126] train-logloss:0.436635+0.003408 test-logloss:0.445521+0.032134
## [127] train-logloss:0.436621+0.003406 test-logloss:0.445497+0.032127
## [128] train-logloss:0.436596+0.003400 test-logloss:0.445468+0.032123
## [129] train-logloss:0.436582+0.003388 test-logloss:0.445436+0.032112
## [130] train-logloss:0.436576+0.003392 test-logloss:0.445434+0.032093
## [131] train-logloss:0.436567+0.003393 test-logloss:0.445440+0.032075
## [132] train-logloss:0.436547+0.003417 test-logloss:0.445441+0.032072
## [133] train-logloss:0.436535+0.003422 test-logloss:0.445450+0.032086
## [134] train-logloss:0.436519+0.003412 test-logloss:0.445430+0.032107
## [135] train-logloss:0.436504+0.003399 test-logloss:0.445392+0.032124
## [136] train-logloss:0.436487+0.003381 test-logloss:0.445357+0.032137
## [137] train-logloss:0.436466+0.003373 test-logloss:0.445326+0.032129
## [138] train-logloss:0.436456+0.003370 test-logloss:0.445339+0.032107
## [139] train-logloss:0.436436+0.003395 test-logloss:0.445325+0.032091
## [140] train-logloss:0.436433+0.003392 test-logloss:0.445337+0.032083
## [141] train-logloss:0.436423+0.003395 test-logloss:0.445332+0.032087
## [142] train-logloss:0.436415+0.003398 test-logloss:0.445337+0.032081
## [143] train-logloss:0.436408+0.003398 test-logloss:0.445337+0.032067
## [144] train-logloss:0.436402+0.003398 test-logloss:0.445333+0.032058
## [145] train-logloss:0.436397+0.003403 test-logloss:0.445333+0.032062
## [146] train-logloss:0.436391+0.003399 test-logloss:0.445339+0.032059
## [147] train-logloss:0.436390+0.003399 test-logloss:0.445334+0.032064
## [148] train-logloss:0.436377+0.003414 test-logloss:0.445332+0.032060
## [149] train-logloss:0.436363+0.003415 test-logloss:0.445320+0.032065
## [150] train-logloss:0.436354+0.003423 test-logloss:0.445305+0.032034
## [151] train-logloss:0.436350+0.003422 test-logloss:0.445306+0.032034
## [152] train-logloss:0.436341+0.003421 test-logloss:0.445329+0.032026
## [153] train-logloss:0.436338+0.003421 test-logloss:0.445321+0.032028
## [154] train-logloss:0.436330+0.003417 test-logloss:0.445336+0.032030
## [155] train-logloss:0.436316+0.003419 test-logloss:0.445311+0.032028
## [156] train-logloss:0.436295+0.003430 test-logloss:0.445293+0.032019
## [157] train-logloss:0.436293+0.003428 test-logloss:0.445296+0.032015

```


[illegible]

[illegible]

[illegible]

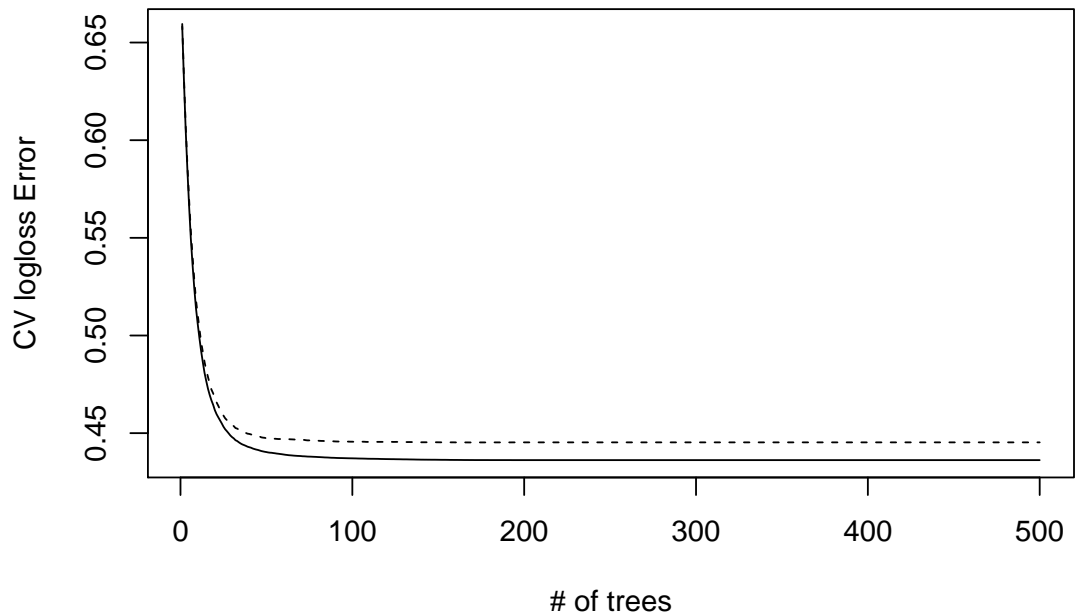
[illegible]

[illegible]

[illegible]

```
## [482] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [483] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [484] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [485] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [486] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [487] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [488] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [489] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [490] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [491] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [492] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [493] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [494] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [495] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [496] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [497] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [498] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
## [499] train-logloss:0.436198+0.003501 test-logloss:0.445282+0.031963
```

```
plot(data.frame(xgb.model.three)[,1], type='l', col='black', ylab='CV logloss Error', xlab=
lines(data.frame(xgb.model.three)[,3], type='l', lty=2, col='black'))
```



As we can see, the minimum of the second model's CV logloss error is $0.434847+0.005887$ which is greater than the first model, $0.434693+0.008852$. The 2nd XGBoost model is a boosted tree model compared to the 1st XGBoost model which is boosted form of logistic regression. 2nd model is getting the minimum of the CV logless error at closer to 100 trees compared to 200 trees for the 1st model, we can say that 2nd model is preferred based on model complexity. Even though the minimum of the CV logless error is slightly higher in the 2nd model as noted earlier, it is negligible in this case and hence the 2nd model is preferred.