



## Individual Coursework Submission Form

### Specialist Masters Programme

<b>Surname: Haixiang</b>	<b>First Name: Yan</b>
<b>MSc in: Quantitative Finance</b>	<b>Student ID number: 190054621</b>
<b>Module Code: SMM636</b>	
<b>Module Title: Machine Learning</b>	
<b>Lecturer: Zhu, Rui</b>	<b>Submission Date: 03 April 2020</b>
<b>Declaration:</b> By submitting this work, I declare that this work is entirely my own except those parts duly identified and referenced in my submission. It complies with any specified word limits and the requirements and regulations detailed in the coursework instructions and any other relevant programme and module documentation. In submitting this work I acknowledge that I have read and understood the regulations and code regarding academic misconduct, including that relating to plagiarism, as specified in the Programme Handbook. I also acknowledge that this work will be subject to a variety of checks for academic misconduct.  We acknowledge that work submitted late without a granted extension will be subject to penalties, as outlined in the Programme Handbook. Penalties will be applied for a maximum of five days lateness, after which a mark of zero will be awarded.	
<b>Marker's Comments (if not being marked on-line):</b>	

**Deduction for Late Submission:**

**Final Mark:**

 %



# **SMM636 Machine Learning (PRD2 A 2019/20)**

## **Individual Assignment**

Yan Haixiang  
March 2020

### Table of Contents

<b>Introduction.....</b>	<b>3</b>
<b>Support Vector Machine (SVM) and ROC Plots .....</b>	<b>2</b>
<b>Test error rates with different methods and tuned parameters .....</b>	<b>3</b>
SVM (kernel = liner) .....	3
SVM (kernel = radial) .....	3
SVM (kernel = polynomial, degree = 2) .....	3
<b>ROC plots and interpretations .....</b>	<b>4</b>
<b>Decision Tree, Random Forest, and ROC Plots .....</b>	<b>5</b>
<b>Decision tree .....</b>	<b>5</b>
Cross-validation (10-fold, repeat 3 times) .....	5
Tree pruning with terminal nodes = 5 .....	5
Additional findings – optimal tree size .....	6
<b>Random Forest .....</b>	<b>5</b>
Test error rate .....	6
Variable Importance .....	7
<b>ROC plots and interpretations .....</b>	<b>7</b>
<b>Conclusion .....</b>	<b>8</b>
<b>Appendix .....</b>	<b>8</b>

## 1. Introduction

This ML assignment explores the robust applications of Support Vector Machine, Decision Tree, as well as Random Forest with the OJ data downloaded from ISLR package. Under each of the three sections, the model is trained with data representing 70% of the entire dataset, while the remaining 30% is reserved for model testing. For each model trained, a test error rate is computed as an indicator of the reliability of that particular model. Besides, as you will see below, with different assumptions of parameters and training methods (especially evident in SVM where three methods, namely linear, radial, and polynomial) applied, the final results differ. The results are plotted for better visualization.

## 2. Support Vector Machine (SVM) and ROC Plots

### 2.1 Test error rates with different methods and tuned parameter

#### 2.1.1 SVM (kernel = linear)

For different cost values tuned, which in this case cost is set to 0.01, 0.10, 1.00, and 10.00 respectively, different error rate results. As shown below, when cost is set to 10.00, we have the lowest error rate (0.1680).

```
- Detailed performance results:
cost      error dispersion
1  0.01  0.1746667  0.05162065
2  0.10  0.1733333  0.04532462
3  1.00  0.1733333  0.04988877
4 10.00  0.1680000  0.04836078
```

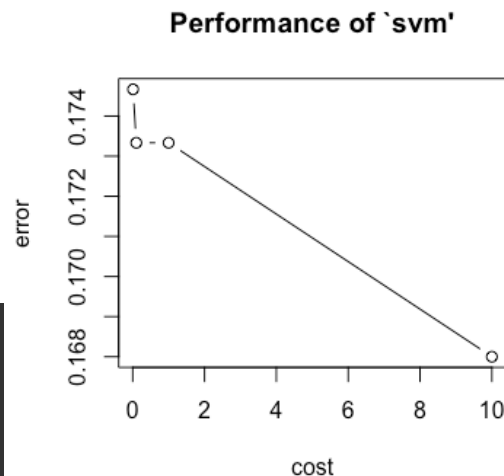


Figure 1, 2 – Test error rates under linear SVM model

#### 2.1.2 SVM (kernel = radial)

Radial is mostly used kernel for SVM. As shown below, when cost is set to 1.00, we have the lowest error rate (0.1853).

```
- Detailed performance results:
cost      error dispersion
1  0.01  0.3893333  0.05989291
2  0.10  0.2666667  0.06942444
3  1.00  0.1853333  0.04283762
4 10.00  0.2040000  0.05931294
```

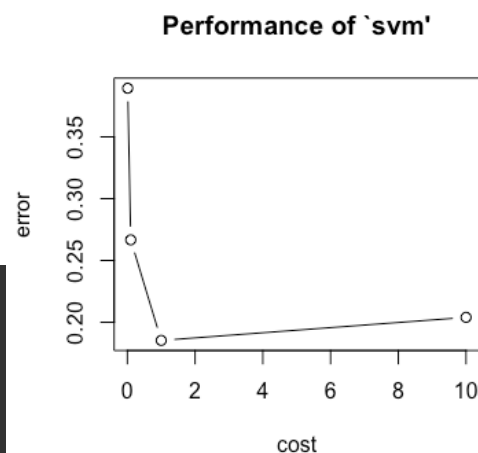


Figure 3, 4 – Test error rates under radial SVM model

## 2.1.3 SVM (kernel = polynomial, degree = 2)

Under this kernel, when cost is set to 10.00, we have the lowest error rate (0.1733).

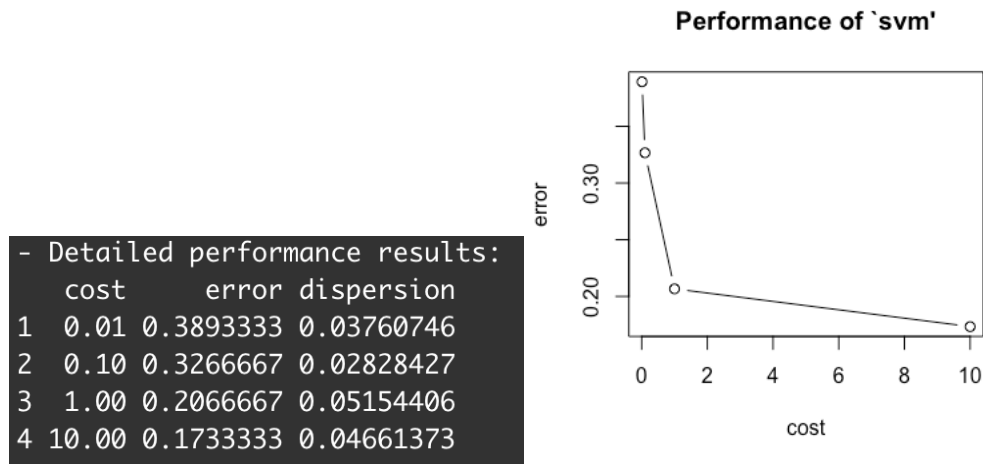


Figure 5, 6 – Test error rates under radial SVM model

In summary, although different kernel models have different optimal cost value for which its error rate is lowest, the general pattern observed is that when cost value increases (at least until 10.00), the error rate decreases.

Besides, it is worth noticing that for each time the model is run, the error rate results are likely to be different. One solution to get a reliable single value with great predictability capacity is to increase the number of times the models are run. Due to time limit of this assignment, this is not further elaborated.

## 2.2 ROC plots and interpretations

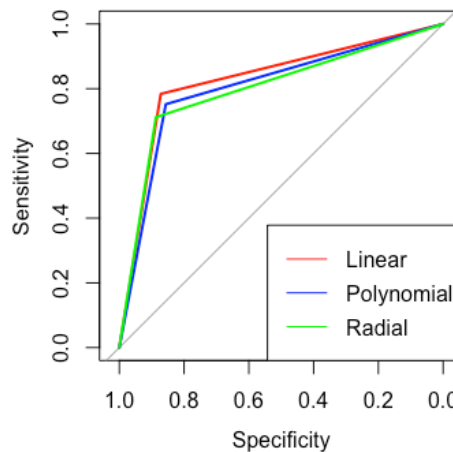


Figure 7 – ROC plotting for three models under SVM

Receiver operating characteristics (ROC) curve serves as a graphical representation of sensitivity and specificity with different classification thresholds. As shown above, the three methods have similar ROC plots and thus areas under curve (AUC), with the liner model slightly better than the other two. Lastly, for all three models, the performance is up to standard, as the AUC is significant.

Besides, non-linear models (such as radial and polynomial) are more likely to result in overfitting problems as they tend to interpret the existing data well but when comes to prediction capabilities, they are disadvantaged.

### 3. Decision Tree, Random Forest, and ROC Plots

#### 3.1 Decision tree

##### 3.1.1 Cross validation (10-fold, repeat 3 times)

Before apply cross-validation to increase prediction accuracy, the number of terminal nodes is 9. *However*, when applying a 10-fold, 3-time cross validation on the tree model, the optimal tree size is reduced to 4 (*Appendix A*). The accuracy rate with different cp values are present below:

```
CART

750 samples
17 predictor
2 classes: 'CH', 'MM'

No pre-processing
Resampling: Cross-Validated (10 fold, repeated 3 times)
Summary of sample sizes: 675, 675, 675, 675, 675, 675, ...
Resampling results across tuning parameters:

cp          Accuracy   Kappa
0.01598174  0.7964753  0.5741827
0.03253425  0.7902233  0.5586241
0.48287671  0.6986713  0.2906605

Accuracy was used to select the optimal model using the largest value.
The final value used for the model was cp = 0.01598174.
```

Figure 8 – Cross-validation Accuracy Rates

The optimal cp value is 0.01598174, with an accuracy rate of 0.7964753.

##### 3.1.2 Tree pruning with terminal node of 5

Despite the results from cross-validation above, I further prune the number of terminal nodes to 5, here is the shape of the pruned tree and corresponding accuracy rates.

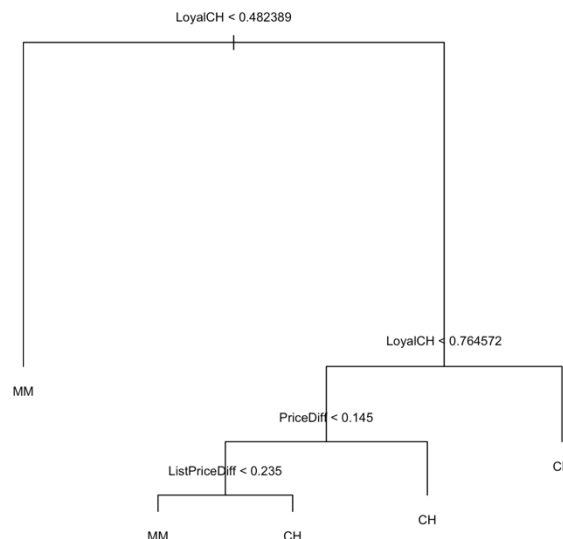


Figure 9 – Tree pruning with number of terminal nodes set to 5

The test error rate of the pruned tree is 0.184375, which is a slight improvement over the cross-validation method above. However, this error rate may not be the

lowest for all the possible tree size we can choose. The optimal number of terminal nodes, though not the goal of this assignment, is further explored below.

### 3.1.3 Additional findings – optimal tree size

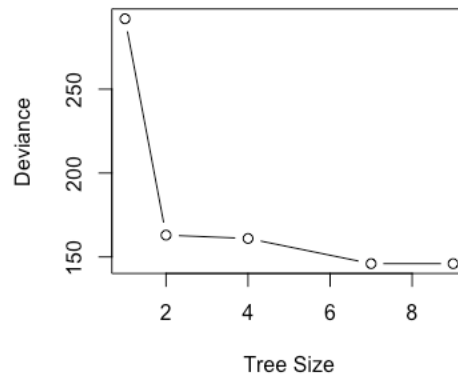


Figure 10 – Optimal Tree Size

The deviation data is set as the indicator to choose the optimal tree size, As shown The deviation value is selected as the parameter to choose the optimal tree size. As shown above, when tree size is 7 or 9, the deviation is at its lowest, which is around 150. Thus a better selection of pruned tree should have a terminal nodes of size 7 or 9, suppose all else are equal.

## 3.2 Random forest

### 3.2.1 Test error rate

Under this section, the parameter “mtry” is set to 1, 2, 3, 4, 5, 6, each with a different accuracy rate, which are illustrated below:

```
Random Forest

750 samples
17 predictor
2 classes: 'CH', 'MM'

No pre-processing
Resampling: Bootstrapped (25 reps)
Summary of sample sizes: 750, 750, 750, 750, 750, 750, ...
Resampling results across tuning parameters:

mtry  Accuracy  Kappa
1     0.7522585  0.4475725
2     0.7943541  0.5588451
3     0.8004569  0.5747206
4     0.7999495  0.5754432
5     0.7987759  0.5729798
6     0.7970948  0.5702088

Accuracy was used to select the optimal model using the largest value.
The final value used for the model was mtry = 3.
```

Figure 11 – Random forest test error rate

When mtry = 3, the accuracy is the highest at 0.8004589. It is worth noticing that for each time the model is run, the results are likely to be different. For this particular run, so long the mtry value is within 2 to 6, the accuracy rates show similar patterns.

### 3.2.2 Variable importance

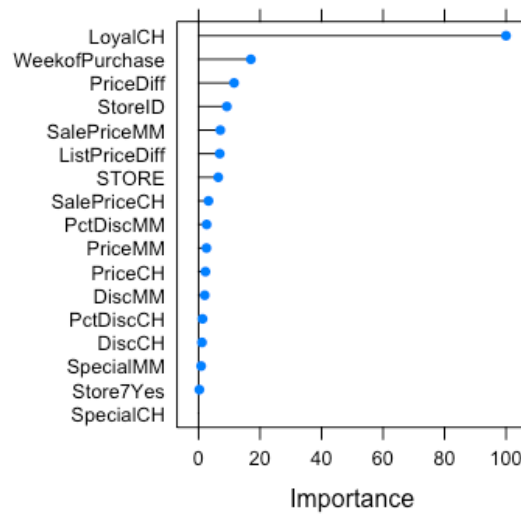


Figure 12 – Variable Importance

Based upon the result that when  $mtry = 3$  the accuracy rate is the best, I plot the variable importance graph to present the variables that influence the prediction outcome the most. As shown in the above graph, the variable “Loyal CH” dominates the chart as the best predicting variable, followed by “Week of Purchase”, “Price Difference”, and “Store ID”, which are of the second tier in terms of importance.

### 3.3 ROC plots and interpretations

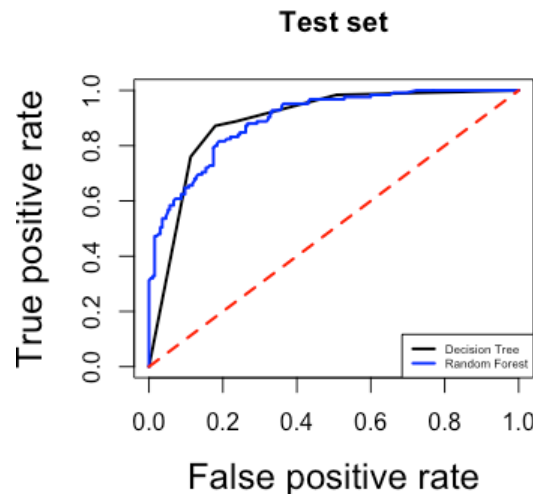


Figure 13 – ROC for decision tree and random forest

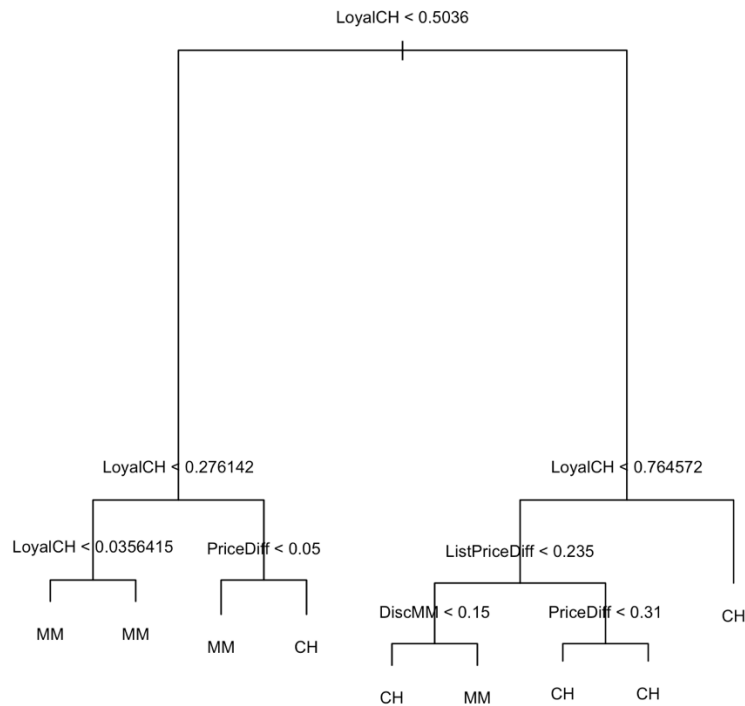
As shown above, both models perform well in this set of data, with significant AUC value, significantly above the benchmark line (red-dotted line), with the Decision Tree model has a AUC of 0.8896821 and Random Forest model has a value of 0.8948923.

## 4. Conclusion

All of the above models have a satisfying prediction capability with accuracy rate concentrating within 0.77 to 0.82. However, this results may be impacted by the particular set of data chosen. For further research, such limitations could be addressed.

## 5. *Appendix*

### Appendix A – Decision Tree without pruning



```

Classification tree:
tree(formula = Purchase ~ ., data = train)
Variables actually used in tree construction:
[1] "LoyalCH"      "PriceDiff"    "ListPriceDiff" "DiscMM"
Number of terminal nodes: 9
Residual mean deviance: 0.7418 = 549.6 / 741
Misclassification error rate: 0.164 = 123 / 750
  
```

### Appendix B – Coding Scripts



```

1  ### R INDIVIDUAL COURSEWORK SVM+DTRF_Haixiang
2  ## Question 1
3  # install packages
4  install.packages("ISLR")
5  install.packages("e1071")
6  install.packages("caret")
7  install.packages("pROC")
8  install.packages("ROCR")
9  library(ISLR)
10 library(caret)
11 library(e1071)
12 library(tree)
13 library(pROC)
14 library(randomForest)
15 library(ROCR)
16
17
18 # load data
19 data(OJ)
20 attach(OJ)
21
22 # data tidiness
23 OJ[complete.cases(OJ), ] #check if there is any missing data
24 OJ
25
26 # split data with 70% training and 30% test data
27 trainIndex <- createDataPartition(OJ$Purchase, p=0.7,list=FALSE, times = 1)
28
29 train<-OJ[trainIndex,]
30 test<-OJ[-trainIndex,]
31 trainlabel<-OJ[trainIndex]
32 testlabel<-OJ[-trainIndex]
33
34 # specify the cost vector
35 set.seed(1234)
36 cost<-c(0.01,0.1,1,10)
37
38 #####
39 #(1) fit support vector classifier
40 set.seed(122)
41 svm_linear = tune(svm,Purchase~.,data=train,
42                 ranges=list(cost=cost),kernel='linear')
43
44 # summary of accuracy rates with corresponding cost values
45 summary(svm_linear)
46 # plot
47 plot(svm_linear)
48
49 #####
50 #(2) fit support vector machine with radial kernel
51 svm_radial = tune(svm,Purchase~.,data=train,

```

```

52         ranges=list(cost=cost),kernel='radial',gamma=0.25)
53
54     # summary of accuracy rates with corresponding cost values
55     summary(svm_radial)
56     #plot
57     plot(svm_radial)
58
59     #####
60     #(3) fit support vector machine with polynomial kernel
61     svm_poly = tune(svm,Purchase~.,data=train,
62                     ranges=list(cost=cost),kernel='polynomial',degree=2)    #specify cost and degree parameters
63
64     # summary of accuracy rates with corresponding cost values
65     summary(svm_poly)
66     #plot
67     plot(svm_poly)
68
69     #####
70     #(4) compute test prediction and accuracy rate for each of the above SVM methods
71     pred1 = predict(svm_linear$best.model,test,decision.values=TRUE)
72     testrate1 = mean(test$Purchase != pred1)
73     pred2 = predict(svm_radial$best.model,test,decision.values=TRUE)
74     testrate2 = mean(test$Purchase != pred2)
75     pred3 = predict(svm_poly$best.model,test,decision.values=TRUE)
76     testrate3 = mean(test$Purchase != pred3)
77
78     # ROC function
79     rocplot =function (pred , truth , order,...){
80         predob = prediction (pred , truth,label.ordering = order)
81         perf = performance (predob,"tpr", "fpr")
82         plot(perf,...)}
83
84     #source("rocplot.R")
85
86     # ROC plotting
87     roc_linear <- roc(response = test$Purchase, predictor =as.numeric(pred1))
88     roc_radial <- roc(response = test$Purchase, predictor =as.numeric(pred2))
89     roc_poly <- roc(response = test$Purchase, predictor =as.numeric(pred3))
90
91     plot(roc_linear,col = c("red"))
92     plot(roc_radial,col = c("blue"), add=TRUE)
93     plot(roc_poly,col = c("green"), add=TRUE)
94
95     # set up legend details
96     legend("bottomright", legend = c("Linear","Polynomial","Radial"), lty = c(1), col = c("red","blue","green"))
97
98     #####
99     #(5) decision tree
100     # install.packages
101     install.packages("tree")

```

```

102 library("tree")
103 install.packages("randomForest")
104 library("randomForest")
105 install.packages("gbm")
106 library("gbm")
107 install.packages("rpart")
108 library("rpart")
109 install.packages("e1071")
110 library("e1071")
111 install.packages("rpart.plot")
112 library("rpart.plot")
113 install.packages("AUC")
114 library(AUC)
115
116 # (5.1) train the model without cross validation or pruning
117 set.seed(13243)
118 OJ_tree = tree(Purchase ~., train)
119 summary(OJ_tree)
120
121 # plot the tree without cross validation or pruning
122 plot(OJ_tree)
123 text(OJ_tree,pretty=1,cex=0.8)
124
125 # evaluate the error rate
126 testrate_dt = mean(predict(OJ_tree, test, type='class')!=test$Purchase)
127 testrate_dt
128
129 # (5.2) train the model with cross-validation (10-fold, repeat 3 times)
130 fitcontrol=trainControl(method = "repeatedcv",
131                          number = 10,
132                          repeats = 3)
133 set.seed(1)
134 OJ_tree_cv=train(train[,-1],
135                  train[,1],
136                  method = "rpart",
137                  trControl = fitcontrol)
138 OJ_tree_cv
139
140 # see the prediction result
141 pred_dt_cv = predict(OJ_tree_cv, test)
142
143 # plot the tree after cross validation
144 plot(OJ_tree_cv$finalModel)
145 text(OJ_tree_cv$finalModel,pretty=1,cex=.8)
146
147 # (5.3) find the optimal tree size
148 #tree pruning
149 OJ_cv = cv.tree(OJ_tree,FUN=prune.misclass)
150 plot(OJ_cv)
151 text(OJ_cv,pretty=1,cex=0.8)

```

```

152
153 # visualize the best tree size by examining deviation
154 plot(OJ_cv$size, OJ_cv$dev, type = "b", xlab = "Tree Size", ylab = "Deviance")
155
156 # prune the tree according to the best tree size based on above analysis
157 prunedTree_best = prune.misclass(OJ_tree, best = OJ_cv$size[which.min(OJ_cv$dev)])
158 pred_dt_best = predict(prunedTree_best, test)
159 testrate_dt_best = mean(predict(prunedTree_best, test, type='class')!=test$Purchase)
160 testrate_dt_best
161
162 # (5.4) prune the tree with terminal nodes of 5
163 prunedTree_five = prune.misclass(OJ_tree, best = 5)
164 pred_dt_five = predict(prunedTree_five, test)
165 testrate_dt_five = mean(predict(prunedTree_five, test, type='class')!=test$Purchase)
166 testrate_dt_five
167
168 # plot the pruned tree with terminal nodes of 5
169 plot(prunedTree_five)
170 text(prunedTree_five, pretty=1, cex=0.8)
171
172 # (6) random forest
173 # train the model under random forest method
174 set.seed(2)
175 mtryGrid=expand.grid(mtry=c(1,2,3,4,5,6))
176 OJ_rf=train(Purchase~., data=train, method="rf",
177             metric="Accuracy",
178             tuneGrid=mtryGrid)
179 OJ_rf
180
181 # prediction with random forest
182 pred_rf = predict(OJ_rf, test, type='prob')
183 pred_rf
184
185 # variable importance
186 # importance(OJ_rf)
187 varImp(OJ_rf)
188
189 # plot variable importance
190 plot(varImp(OJ_rf))
191
192 # (7) ROC plot for decision tree and random forest
193 dt_roc = rocplot(pred_dt_five[,2], test[,1], order=c("CH", "MM"), col="black", lwd=2, cex.lab=1.5, cex.axis=1.5, main="Test set")
194 rf_roc = rocplot(pred_rf[,2], test[,1], order=c("CH", "MM"), add=TRUE, col="blue", lwd=2, cex.lab=1.5, cex.axis=1.5)
195
196 # AUC values
197 auc(roc(pred_dt_five[,2], test[,1]))
198 auc(roc(pred_rf[,2], test[,1]))
199
200 # set up legend details
201 legend("bottomright",
202       legend = c("Decision Tree", "Random Forest"),
203       col=c("black", "blue"), cex=0.5, lty=1, lwd=2)
204
205 # set the axis scales
206 x=seq(0,1,0.01); y=x
207
208 # add reference line
209 lines(x,y, lwd=2, col="red", lty=2)
210
211
212
213

```