

MACHINE LEARNING AND DATA MINING PROJECT

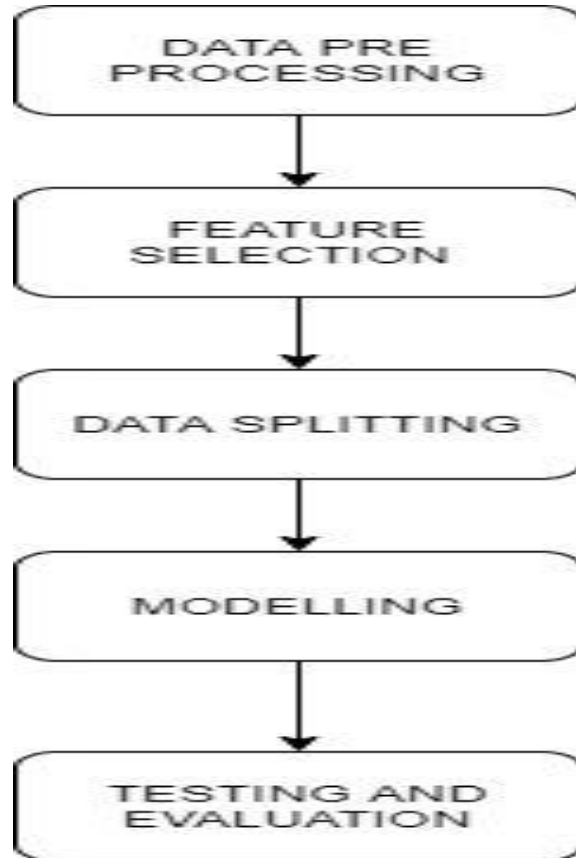
VEHICLE LOAN DEFAULT PREDICTION

REGISTRATION NUMBER	NAME
BL.EN.U4CSE17072	SARATH CHANDRA .M
BL.EN.U4CSE17076	M.M. SHOBANDARI
BL.EN.U4CSE17121	SUNIL CHOWDARY .B

ABSTRACT

- Financial institutions incur significant losses due to the default of vehicle loans. This has led to the tightening up of vehicle loan underwriting and increased vehicle loan rejection rates.
- The need for a better credit risk scoring model is also raised by these institutions. This warrants a study to estimate the determinants of vehicle loan default.
- A financial institution has hired you to accurately predict the probability of loanee/borrower defaulting on a vehicle loan in the first EMI (Equated Monthly Installments) on the due date.
- Doing so will ensure that clients capable of repayment are not rejected and important determinants can be identified which can be further used for minimising the default rates.

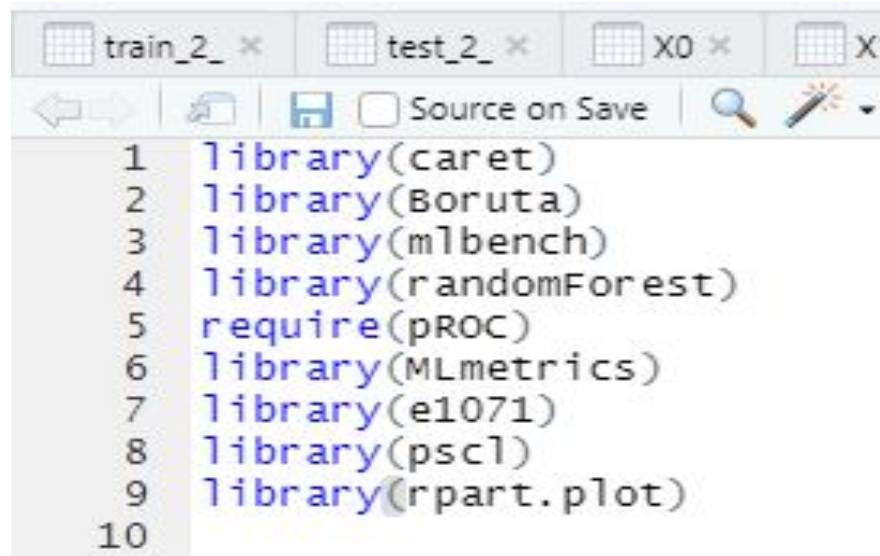
IMPLEMENTATION



MODELS :

- RANDOM FORESTS
- LOGISTIC REGRESSION
- DECISION TREE
- NAIVE BAYES

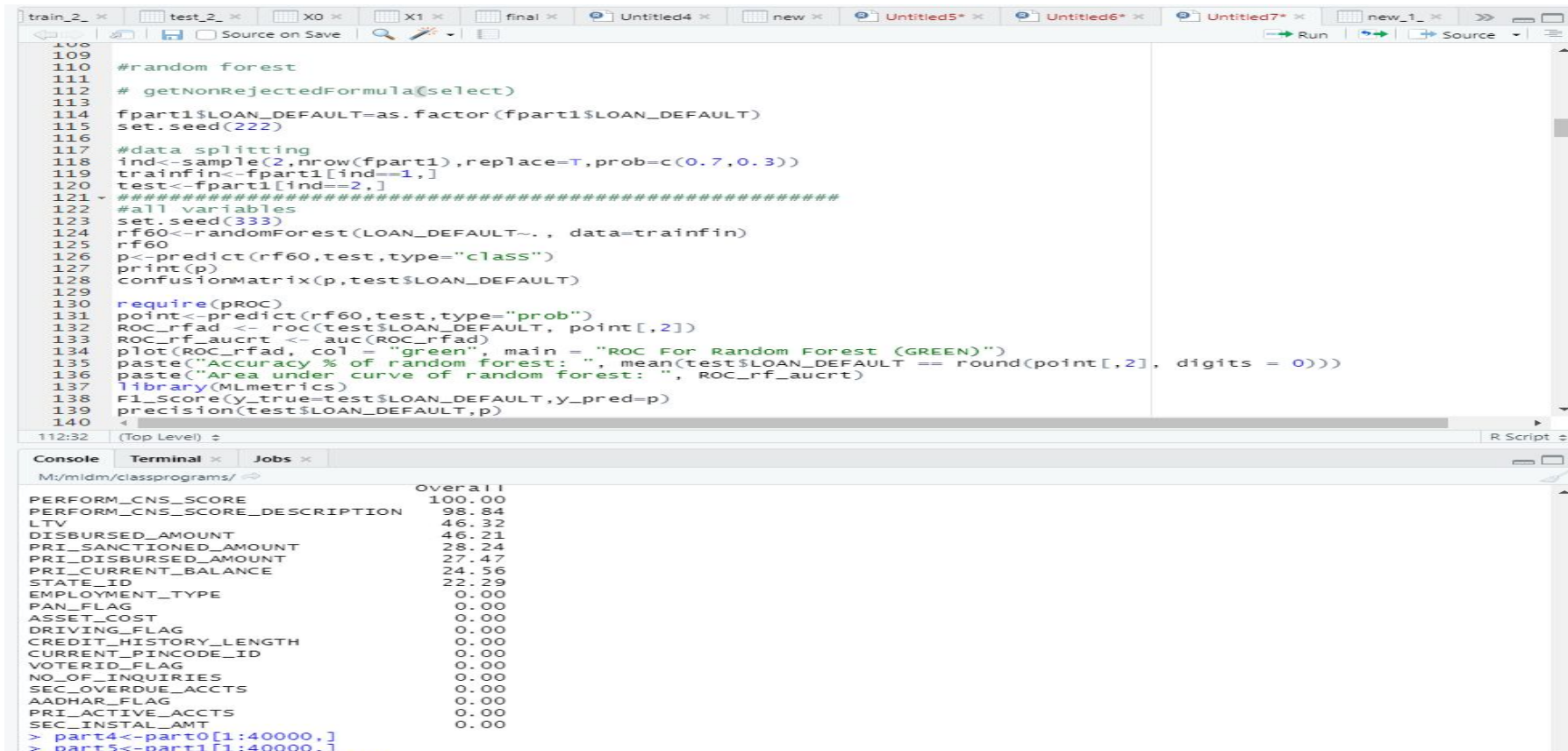
IMPORTED LIBRARIES



A screenshot of the RStudio code editor window. The title bar shows four open files: 'train_2_', 'test_2_', 'X0', and 'X'. The toolbar includes navigation arrows, a source icon, a save icon, a checkbox labeled 'Source on Save', a search icon, and a help icon. The code editor displays the following R code:

```
1 library(caret)
2 library(Boruta)
3 library(mlbench)
4 library(randomForest)
5 require(pROC)
6 library(MLmetrics)
7 library(e1071)
8 library(psc1)
9 library(rpart.plot)
10
```

SELECTING FEATURES USING BORUTA

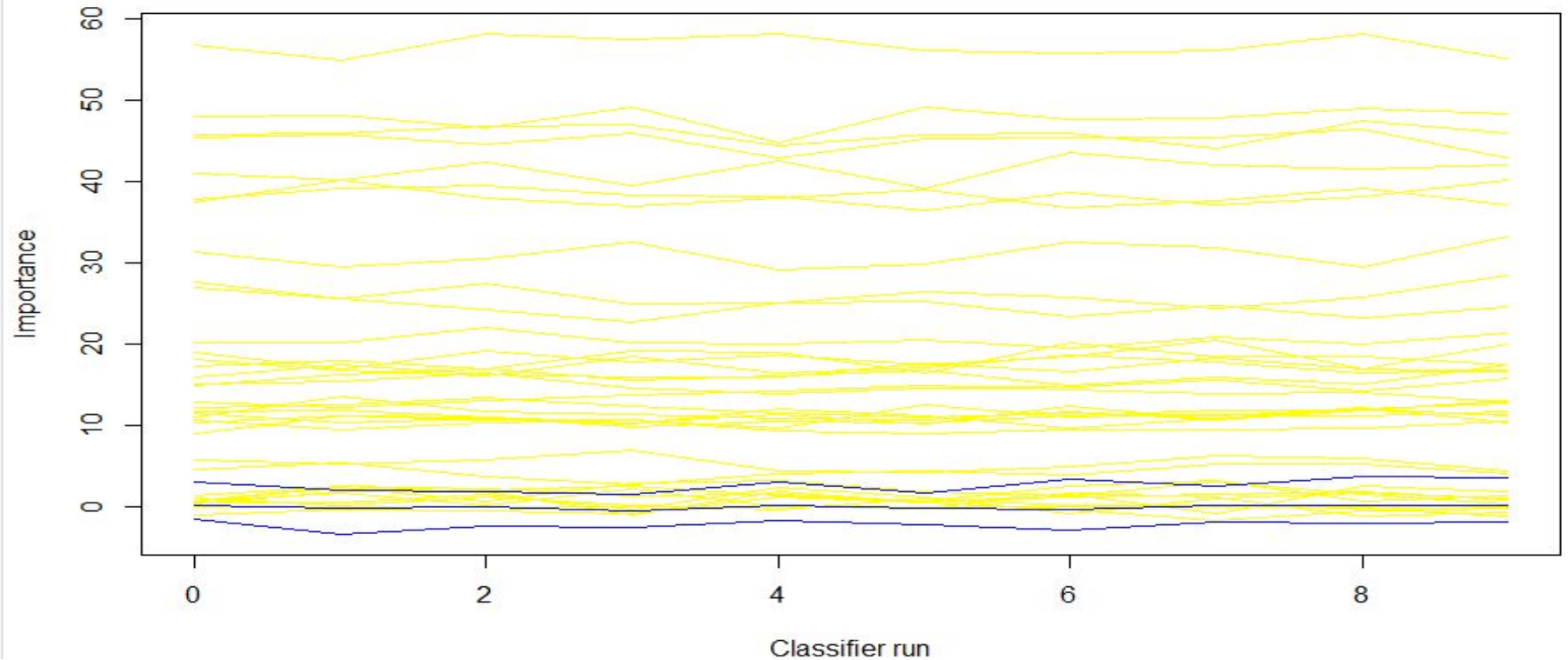


```
108 #random forest
109
110 # getNonRejectedFormula(select)
111
112 fpart1$LOAN_DEFAULT=as.factor(fpart1$LOAN_DEFAULT)
113 set.seed(222)
114
115 #data splitting
116 ind<-sample(2,nrow(fpart1),replace=T,prob=c(0.7,0.3))
117 trainfin<-fpart1[ind==1,]
118 test<-fpart1[ind==2,]
119 #####
120 #all variables
121 set.seed(333)
122 rf60<-randomForest(LOAN_DEFAULT~., data=trainfin)
123 rf60
124 p<-predict(rf60,test,type="class")
125 print(p)
126 confusionMatrix(p,test$LOAN_DEFAULT)
127
128 require(pROC)
129 point<-predict(rf60,test,type="prob")
130 ROC_rfad <- roc(test$LOAN_DEFAULT, point[,2])
131 ROC_rf_aucrt <- auc(ROC_rfad)
132 plot(ROC_rfad, col = "green", main = "ROC For Random Forest (GREEN)")
133 paste("Accuracy % of random forest: ", mean(test$LOAN_DEFAULT == round(point[,2], digits = 0)))
134 paste("Area under curve of random forest: ", ROC_rf_aucrt)
135 library(MLmetrics)
136 F1_Score(y_true=test$LOAN_DEFAULT,y_pred=p)
137 precision(test$LOAN_DEFAULT,p)
138
139 112:32 (Top Level) R Script
```

Console

```
M:/midm/classprograms/
Overall
PERFORM_CNS_SCORE 100.00
PERFORM_CNS_SCORE_DESCRIPTION 98.84
LTV 46.32
DISBURSED_AMOUNT 46.21
PRI_SANCTIONED_AMOUNT 28.24
PRI_DISBURSED_AMOUNT 27.47
PRI_CURRENT_BALANCE 24.56
STATE_ID 22.29
EMPLOYMENT_TYPE 0.00
PAN_FLAG 0.00
ASSET_COST 0.00
DRIVING_FLAG 0.00
CREDIT_HISTORY_LENGTH 0.00
CURRENT_PINCODE_ID 0.00
VOTERID_FLAG 0.00
NO_OF_INQUIRIES 0.00
SEC_OVERDUE_ACCTS 0.00
AADHAR_FLAG 0.00
PRI_ACTIVE_ACCTS 0.00
SEC_INSTAL_AMT 0.00
> part4<-part0[1:40000,]
> part5<-part1[1:40000,]
```

SELECTING FEATURES USING BORUTA

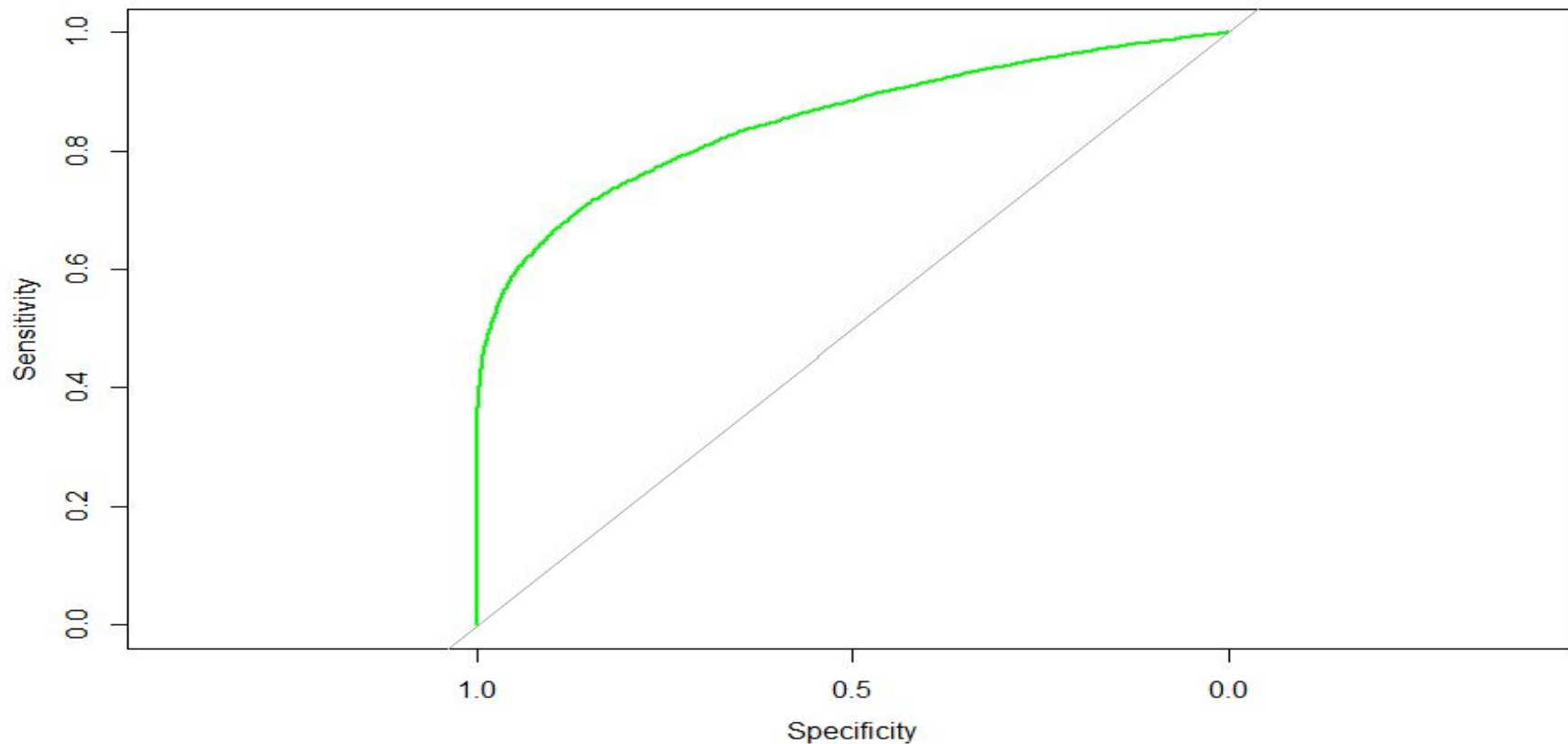


RANDOM FOREST WITH ALL FEATURES

```
train_2_ x  test_2_ x  X0 x  X1 x  final x  Untitled4 x  new x  Untitled5* x  Untitled6* x  Untitled7* x  new_1_ x  >>
Source on Save  Run  Source
126 p<-predict(rf60,test,type="class")
127 print(p)
128 confusionMatrix(p,test$LOAN_DEFAULT)
129
130 require(pROC)
131 point<-predict(rf60,test,type="prob")
132 ROC_rfad <- roc(test$LOAN_DEFAULT, point[,2])
133 ROC_rf_aucrt <- auc(ROC_rfad)
134 plot(ROC_rfad, col = "green", main = "ROC For Random Forest (GREEN)")
135 paste("Accuracy % of random forest: ", mean(test$LOAN_DEFAULT == round(point[,2], digits = 0)))
136 paste("Area under curve of random forest: ", ROC_rf_aucrt)
137 library(MLmetrics)
138 f1_score(y_true=test$LOAN_DEFAULT,y_pred=p)
139 precision(test$LOAN_DEFAULT,p)
140 recall(test$LOAN_DEFAULT,p)
141
142 #####
143 #27 variables
144 rf27<-randomForest(LOAN_DEFAULT~ DISBURSED_AMOUNT + ASSET_COST + PERFORM_CNS_SCORE +
145 PRI_CURRENT_BALANCE + PRI_SANCTIONED_AMOUNT + PRI_DISBURSED_AMOUNT +
146 PRIMARY_INSTALL_AMT + LTV + BRANCH_ID + SUPPLIER_ID + MANUFACTURER_ID +
147 CURRENT_PINCODE_ID + EMPLOYMENT_TYPE + STATE_ID + EMPLOYEE_CODE_ID +
148 AADHAR_FLAG + PAN_FLAG + VOTERID_FLAG + DRIVING_FLAG + PERFORM_CNS_SCORE_DESCRIPTION +
149 PRI_NO_OF_ACCTS + PRI_ACTIVE_ACCTS + PRI_OVERDUE_ACCTS +
150 NEW_ACCTS_IN_LAST_SIX_MONTHS + AVERAGE_ACCT_AGE + CREDIT_HISTORY_LENGTH +
151 NO_OF_INQUIRIES , data=trainfin)
152
153 rf27
154 library(e1071)
155
156 #testing
157 p1<-predict(rf27.test,type="class")
158
130:1 (Untitled) s  R Script s
Console  Terminal x  Jobs x
M:\midm\class\programs\
Detection Rate : 0.4359
Detection Prevalence : 0.5911
Balanced Accuracy : 0.7818
'Positive' Class : 0
> require(pROC)
> point<-predict(rf60,test,type="prob")
> ROC_rfad <- roc(test$LOAN_DEFAULT, point[,2])
> setting levels: control = 0, case = 1
> setting direction: controls < cases
> ROC_rf_aucrt <- auc(ROC_rfad)
> plot(ROC_rfad, col = "green", main = "ROC For Random Forest (GREEN)")
> paste("Accuracy % of random forest: ", mean(test$LOAN_DEFAULT == round(point[,2], digits = 0)))
[1] "Accuracy % of random forest: 0.781572809063645"
> paste("Area under curve of random forest: ", ROC_rf_aucrt)
[1] "Area under curve of random forest: 0.85145596904385"
> library(MLmetrics)
> f1_score(y_true=test$LOAN_DEFAULT,y_pred=p)
[1] 0.7997096
> precision(test$LOAN_DEFAULT,p)
[1] 0.8735392
> recall(test$LOAN_DEFAULT,p)
[1] 0.7373873
>
```


RANDOM FOREST WITH ALL FEATURES

ROC For Random Forest (GREEN)



RANDOM FOREST WITH 27 SELECTED FEATURES

```
151 rf27
152 NO_OF_INQUIRIES , data=trainin)
153
154 library(e1071)
155
156 #testing
157 p1<-predict(rf27,test,type="class")
158 confusionMatrix(p1,test$LOAN_DEFAULT)
159
160 require(pROC)
161 p3<-predict(rf27,test,type="prob")
162 ROC_rf <- roc(test$LOAN_DEFAULT, p3[,2])
163 ROC_rf_auc <- auc(ROC_rf)
164 plot(ROC_rf, col = "green", main = "ROC For Random Forest (GREEN)")
165 paste("Accuracy % of random forest: ", mean(test$LOAN_DEFAULT == round(p3[,2], digits = 0)))
166 paste("Area under curve of random forest: ", ROC_rf_auc)
167 library(MLmetrics)
168 F1_Score(y_true=test$LOAN_DEFAULT,y_pred=p1)
169 precision(test$LOAN_DEFAULT,p1)
170 recall(test$LOAN_DEFAULT,p1)
171
172 #####
173 ##20 variables
174 rf20<-randomForest(LOAN_DEFAULT~PERFORM_CNS_SCORE +
175                     PERFORM_CNS_SCORE_DESCRIPTION +
176                     LTV +
177                     DISBURSED_AMOUNT +
178                     PRI_SANCTIONED_AMOUNT +
179                     PRI_DISBURSED_AMOUNT +
180                     PRI_CURRENT_BALANCE +
181                     STATE_ID +
182                     EMPLOYMENT_TYPE +
183
159:1 (Untitled) R Script :
```

Console

```
M:\mldm\classprograms/
> confusionMatrix(p1,test$LOAN_DEFAULT)
Confusion Matrix and Statistics

      Reference
Prediction 0      1
0    10444    3618
1     1536    8410

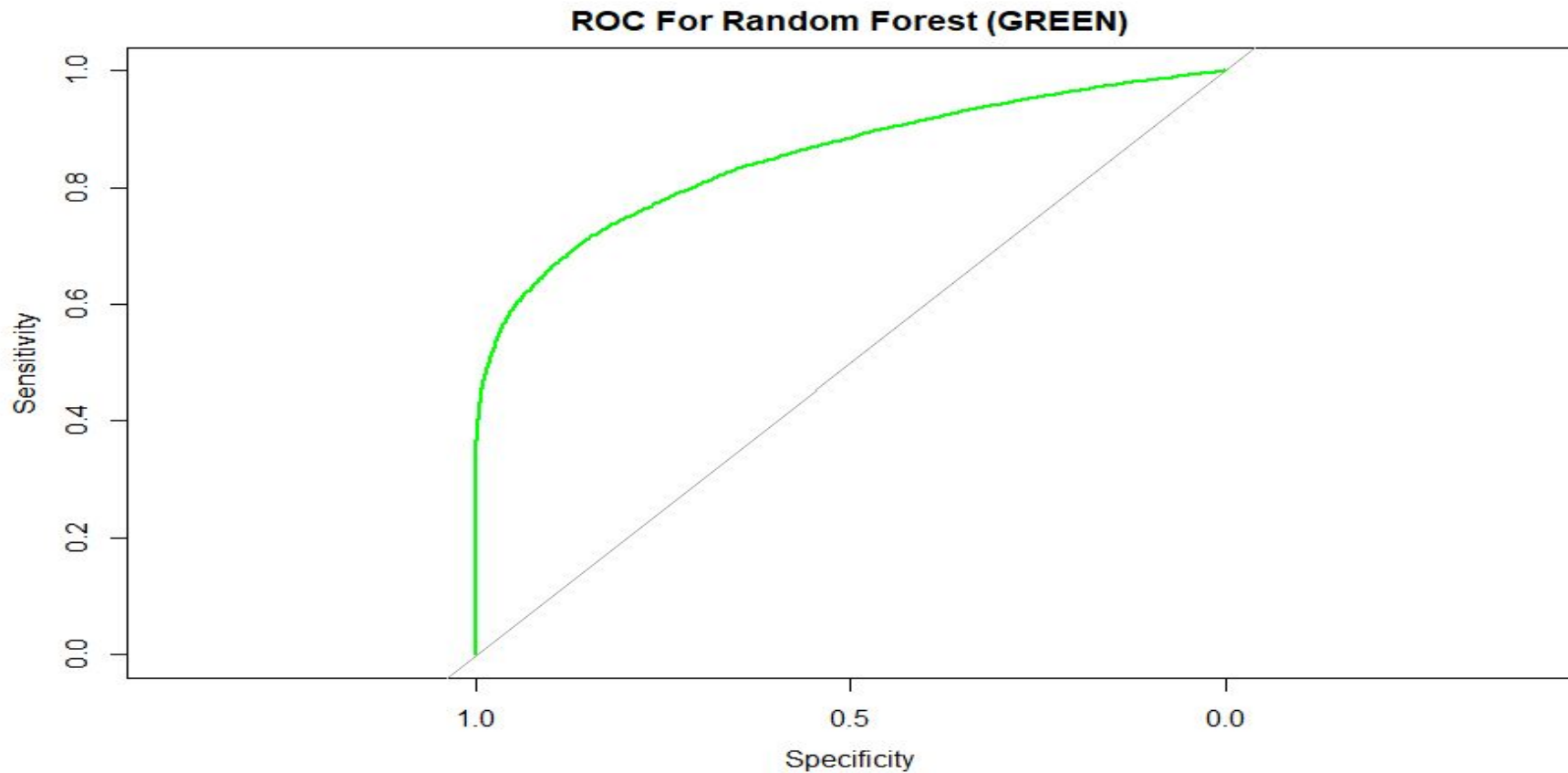
      Accuracy : 0.7853
      95% CI   : (0.7801, 0.7905)
      No Information Rate : 0.501
      P-Value [ACC > NIR] : < 2.2e-16

      Kappa : 0.5708

      Mcnemar's Test P-Value : < 2.2e-16

      Sensitivity : 0.8718
      Specificity : 0.6992
      Pos Pred Value : 0.7427
      Neg Pred Value : 0.8456
      Prevalence : 0.4990
      Detection Rate : 0.4350
      Detection Prevalence : 0.5857
      Balanced Accuracy : 0.7855
```

RANDOM FOREST WITH 27 SELECTED FEATURES



RANDOM FOREST WITH 20 SELECTED FEATURES

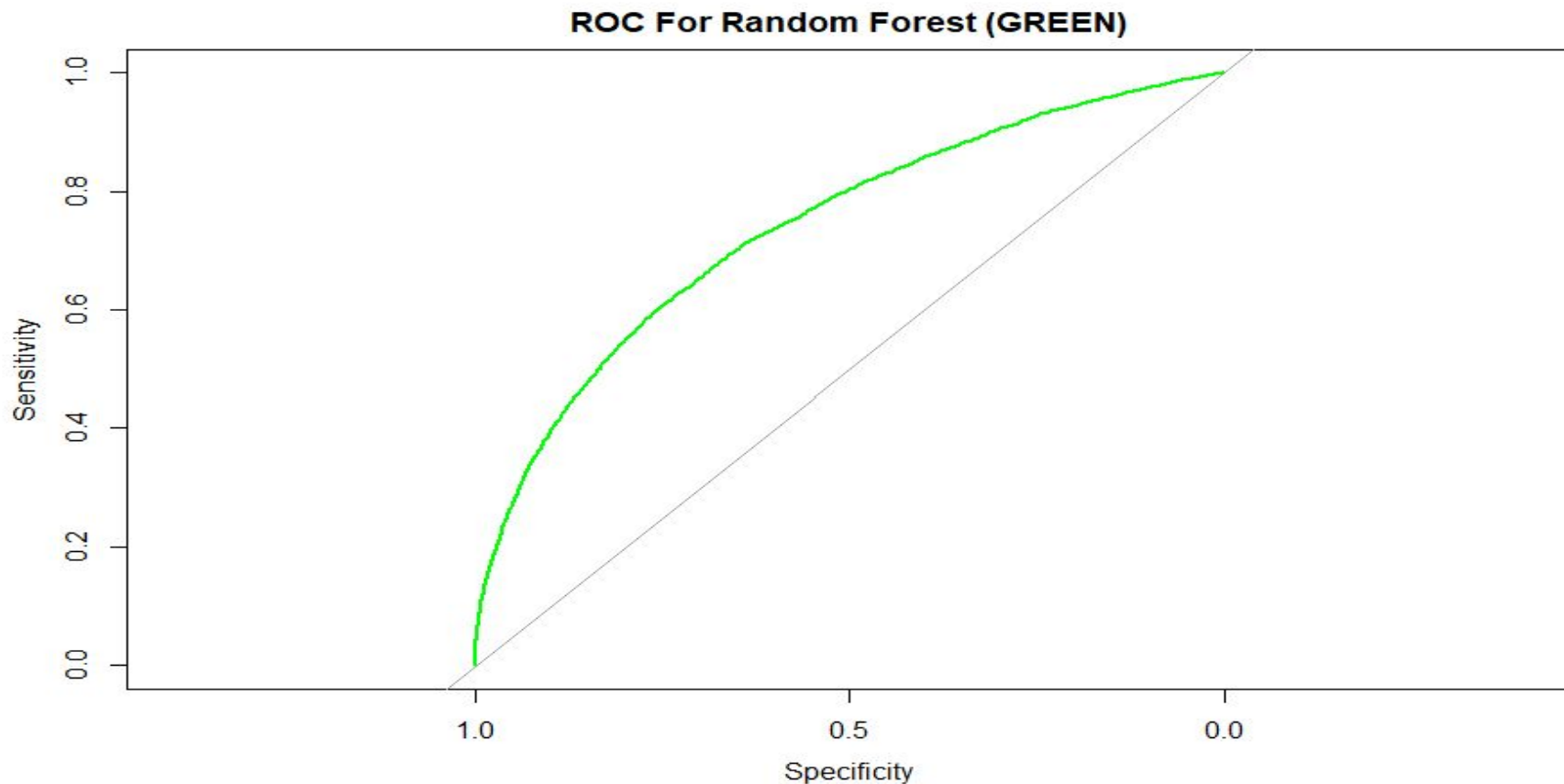
```
train_2_ x  test_2_ x  X0 x  X1 x  final x  Untitled4 x  new x  Untitled5* x  Untitled6* x  Untitled7* x  new_1_ x  >>
Source on Save  Run  Source

192 rf20      PRI_ACTIVE_ACCTS      +
193 SEC_INSTAL_AMT      , data=trainfin)
194
195 #testing
196 p2<-predict(rf20,test,type="class")
197 confusionMatrix(p2,test$LOAN_DEFAULT)
198
199 library(ROCR)
200
201 #metrics
202 require(pROC)
203 point2<-predict(rf20,test,type="prob")
204 ROC_rfas <- roc(test$LOAN_DEFAULT, point2[,2])
205 ROC_rf_aucas <- auc(ROC_rfas)
206 plot(ROC_rfas, col = "green", main = "ROC For Random Forest (GREEN)")
207 paste("Accuracy % of random forest: ", mean(test$LOAN_DEFAULT == round(point2[,2], digits = 0)))
208 paste("Area under curve of random forest: ", ROC_rf_aucas)
209 library(MLmetrics)
210 F1_Score(y_true=test$LOAN_DEFAULT,y_pred=p2)
211 precision(test$LOAN_DEFAULT,p2)
212 recall(test$LOAN_DEFAULT,p2)
213
214
215
216 #####
217
218 #logistic regression
219
220
221
222
223
224 4
225 (Untitled)  R Scrip

Console  Terminal x  Jobs x
M:/midm/classprograms/

> library(ROCR)
warning message:
package 'ROCR' was built under R version 4.0.3
>
> #metrics
> require(pROC)
> point2<-predict(rf20,test,type="prob")
> ROC_rfas <- roc(test$LOAN_DEFAULT, point2[,2])
Setting levels: control = 0, case = 1
Setting direction: controls < cases
> ROC_rf_aucas <- auc(ROC_rfas)
> plot(ROC_rfas, col = "green", main = "ROC For Random Forest (GREEN)")
> paste("Accuracy % of random forest: ", mean(test$LOAN_DEFAULT == round(point2[,2], digits = 0)))
[1] "Accuracy % of random forest: 0.677149283572143"
> paste("Area under curve of random forest: ", ROC_rf_aucas)
[1] "Area under curve of random forest: 0.739740560145415"
> library(MLmetrics)
> F1_Score(y_true=test$LOAN_DEFAULT,y_pred=p2)
[1] 0.6846642
> precision(test$LOAN_DEFAULT,p2)
[1] 0.7028381
> recall(test$LOAN_DEFAULT,p2)
[1] 0.6674065
>
```

RANDOM FOREST WITH 20 SELECTED FEATURES



RESULTS OF ALL IMPLEMENTATIONS

ACCURACY :

MODEL	ALL FEATURES	27 FEATURES	20 FEATURES
RANDOM FOREST	78.17	78.53	67.69
LOGISTIC REGRESSION	66.05	66.05	61.72
DECISION TREE	67.04	67.04	62.45
NAIVE BAYES	64.51	58.16	58.18

RESULTS OF ALL IMPLEMENTATIONS

AUC :

MODEL	ALL FEATURES	27 FEATURES	20 FEATURES
RANDOM FOREST	85.14	85.42	73.97
LOGISTIC REGRESSION	66.05	66.05	61.72
DECISION TREE	69.71	67.90	63.58
NAIVE BAYES	70.30	70.99	62.21

RESULTS OF ALL IMPLEMENTATIONS

F1 SCORE :

MODEL	ALL FEATURES	27 FEATURES	20 FEATURES
RANDOM FOREST	0.7997	0.8020	0.6800
LOGISTIC REGRESSION	0.6679	0.6678	0.6104
DECISION TREE	0.7305	0.7305	0.6091
NAIVE BAYES	0.6494	0.6967	0.6971

REFERENCES

- https://rpubs.com/Monesh_23992/Bank_Loan_Default - Reference model for the whole project
- <https://www.analyticsvidhya.com/blog/2016/03/select-important-variables-boruta-package/> - A description for boruta feature selection
- <https://towardsdatascience.com/random-forest-in-r-f66adf80ec9> - Used for implementing random forest.