

Project: Creditworthiness

Complete each section. When you are ready, save your file as a PDF document and submit it here: <https://classroom.udacity.com/nanodegrees/nd008/parts/11a7bf4c-2b69-47f3-9aec-108ce847f855/project>

Step 1: Business and Data Understanding

Provide an explanation of the key decisions that need to be made. (250 word limit)

Key Decisions:

Answer these questions

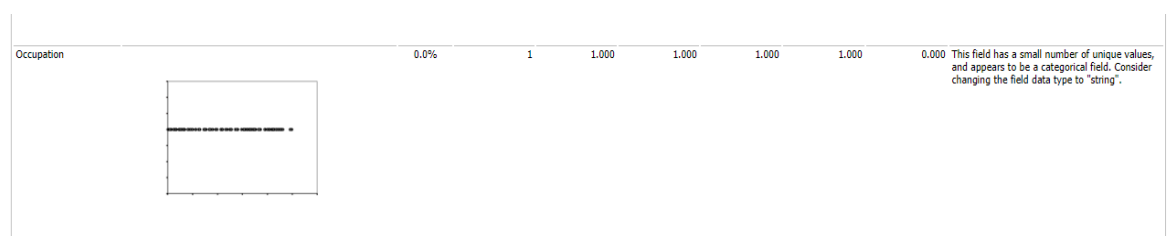
- What decisions needs to be made?
 - ➔ As a business analyst, we need to find out which bank customers are trustworthy to be approved for a loan and which are not.
- What data is needed to inform those decisions?
 - ➔ Past variables of the bank customers for which the bank has either approved or denied the loan ,like age, balance, previous loans, average monthly deposits etc.
 - ➔ Data of the customers , on which we need to score our model and predict the outcomes.
- What kind of model (Continuous, Binary, Non-Binary, Time-Series) do we need to use to help make these decisions?
 - ➔ As there will be only two possible outcomes in the our predictions, so the model we will build will be a binary model.

Step 2: Building the Training Set

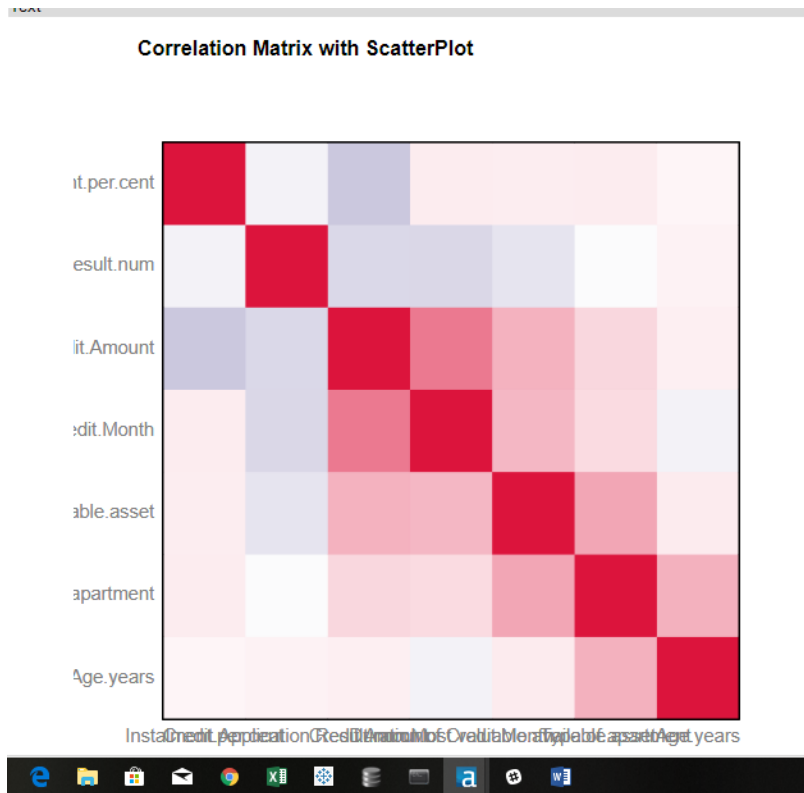
*Build your training set given the data provided to you. The data has been cleaned up for you already so you shouldn't **need to convert any data fields to the appropriate data types.***

Here are some guidelines to help guide your data cleanup:

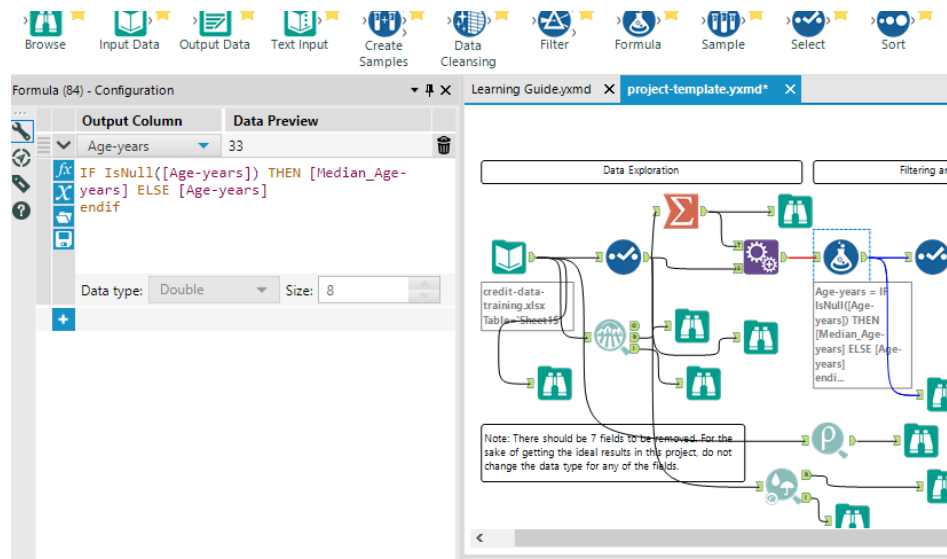
DATA CLEANING



- ➔ From the above info from field summary tool, we can see that concurrent-credits and occupation columns have only one unique value, so it was removed from the data.
- ➔ Guarantors, foreign-worker and no of dependents columns have very low variability, so these columns were also removed.
- ➔ Duration in current address column have missing values more than half of the total rows, so it also needs to be removed.
- ➔ Telephone column is also dropped, due to no logical reason to include the variable.



➔ From the association-analysis tool we can see that no variable is correlated with each other, that is correlation of 0.7 and higher.



➔ Age-years columns have 2% missing values, and these values are filled with the median of Age-years column because In the representation , we can see that data is not normally distributed, it is skewed towards the right.

Step 3: Train your Classification Models

First, create your Estimation and Validation samples where 70% of your dataset should go to Estimation and 30% of your entire dataset should be reserved for Validation. Set the Random Seed to 1.

Create all of the following models: Logistic Regression, Decision Tree, Forest Model, Boosted Model

Answer these questions for **each model** you created:

- Which predictor variables are significant or the most important? Please show the p-values or variable importance charts for all of your predictor variables.
- Validate your model against the Validation set. What was the overall percent accuracy? Show the confusion matrix. Are there any bias seen in the model's predictions?

You should have four sets of questions answered. (500 word limit)

→ Logistic Regression Model :-

- I have used stepwise to automate the logit process.

1 Fields | Records 1 to 10

Report for Logistic Regression Model Stepwise

Basic Summary

Call:
glm(formula = Credit.Application.Result ~ Account.Balance + Payment.Status.of.Previous.Credit + Purpose + Credit.Amount + Length.of.current.employment + Instalment.per.cent + Most.valuable.available.asset, family = binomial(logit), data = the.data)

Deviance Residuals:

	Min	1Q	Median	3Q	Max
	-2.289	-0.713	-0.448	0.722	2.454

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.9621914	6.837e-01	-4.3326	1e-05 ***
Account.BalanceSome Balance	-1.6053228	3.067e-01	-5.2344	1.65e-07 ***
Payment.Status.of.Previous.CreditPaid Up	0.2360857	2.977e-01	0.7930	0.42775
Payment.Status.of.Previous.CreditSome Problems	1.2154514	5.151e-01	2.3595	0.0183 **
PurposeNew car	-1.6993164	6.142e-01	-2.7668	0.00566 ***
PurposeOther	-0.3257637	8.179e-01	-0.3983	0.69042
PurposeUsed car	-0.7645820	4.004e-01	-1.9096	0.05618 .
Credit.Amount	0.0001704	5.733e-05	2.9716	0.00296 ***
Length.of.current.employment4-7 yrs	0.3127022	4.587e-01	0.6817	0.49545
Length.of.current.employment< 1yr	0.8125785	3.874e-01	2.0973	0.03596 **
Instalment.per.cent	0.3016731	1.350e-01	2.2340	0.02549 **
Most.valuable.available.asset	0.2650267	1.425e-01	1.8599	0.06289 .

Significance codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
(Dispersion parameter for binomial taken to be 1)

8 Null deviance: 413.16 on 349 degrees of freedom
Residual deviance: 328.55 on 338 degrees of freedom
McFadden R-Squared: 0.2048, Akaike Information Criterion 352.5

9 Number of Fisher Scoring iterations: 5

10 Type II Analysis of Deviance Tests

- We can see account balance, purpose and credit amount are the most significant predictor variables.

Record Layout

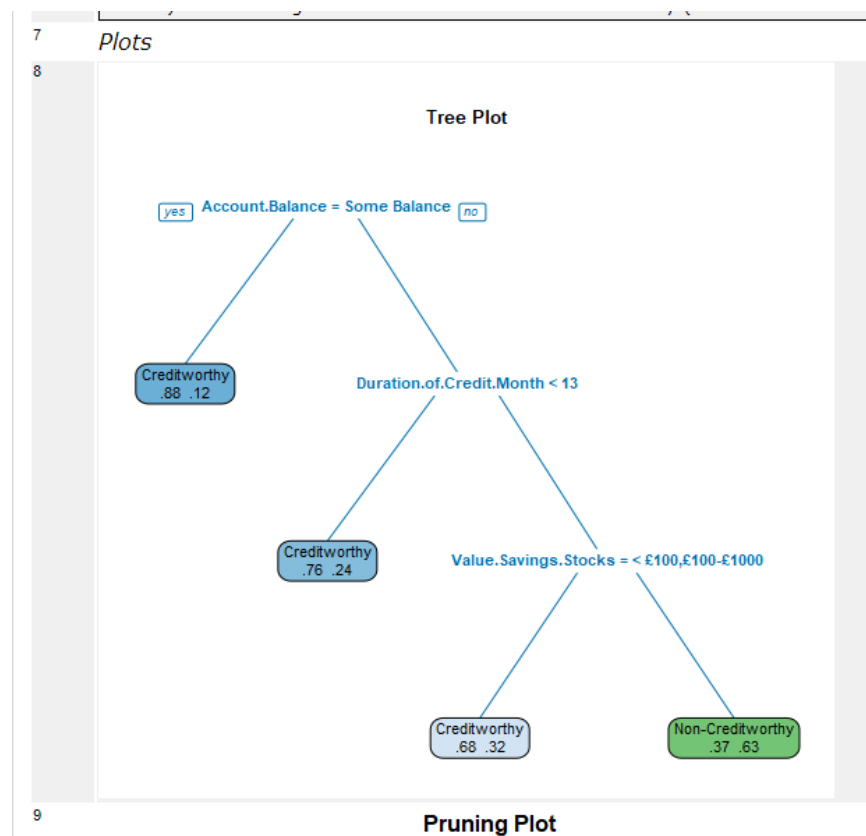
Model Comparison Report						
Fit and error measures						
Model	Accuracy	F1	AUC	Accuracy_Creditworthy	Accuracy_Non-Creditworthy	
Stepwise	0.7600	0.8364	0.7306	0.8762	0.4889	

Model: model names in the current comparison.
Accuracy: overall accuracy, number of correct predictions of all classes divided by total sample number.
Accuracy_[class name]: accuracy of Class [class name] is defined as the number of cases that are **correctly** predicted to be Class [class name] divided by the total number of cases that actually belong to Class [class name], this measure is also known as *recall*.
AUC: area under the ROC curve, only available for two-class classification.
F1: F1 score, $2 * \text{precision} * \text{recall} / (\text{precision} + \text{recall})$. The *precision* measure is the percentage of actual members of a class that were predicted to be in that class divided by the total number of cases predicted to be in that class. In situations where there are three or more classes, average precision and average recall values across classes are used to calculate the F1 score.

Confusion matrix of Stepwise		
	Actual_Creditworthy	Actual_Non-Creditworthy
Predicted_Creditworthy	92	23
Predicted_Non-Creditworthy	13	22

- The overall accuracy of the stepwise model is 76% and we can see that the model is biased for predicting non-creditworthy customers with 48.89% accuracy.

→ Decision Tree



- This is the tree generated by decision tree method.

Summary Report for Decision Tree Model Decision_Tree

Call:

```
rpart(formula = Credit.Application.Result ~ Account.Balance + Duration.of.Credit.Month + Payment.Status.of.Previous.Credit + Purpose + Credit.Amount + Value.Savings.Stocks + Length.of.current.employment + Instalment.per.cent + Most.valuable.available.asset + Age.years + Type.of.apartment + No.of.Credits.at.this.Bank, data = the.data, minsplit = 20, minbucket = 7, usesurrogate = 2, xval = 10, maxdepth = 20, cp = 1e-05)
```

Model Summary

Variables actually used in tree construction:

[1] Account.Balance Duration.of.Credit.Month Value.Savings.Stocks

Root node error: 97/350 = 0.27714

n = 350

Pruning Table

Level	CP	Num Splits	Rel Error	X Error	X Std Dev
1	0.068729	0	1.00000	1.00000	0.086326
2	0.041237	3	0.79381	0.92784	0.084295

Leaf Summary

node), split, n, loss, yval, (yprob)

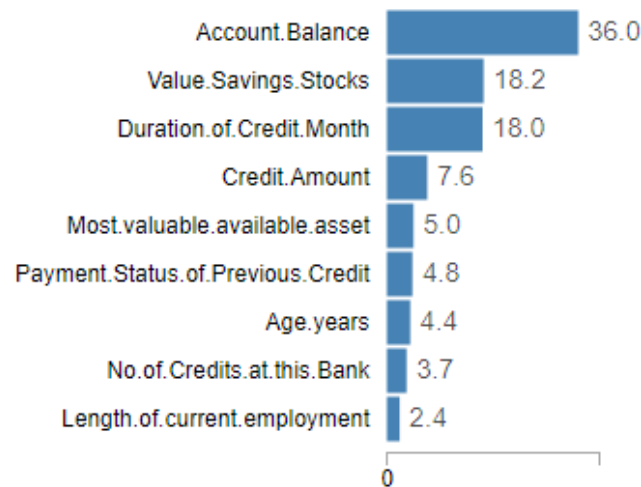
* denotes terminal node

- 1) root 350 97 Creditworthy (0.7228571 0.2771429)
- 2) Account.Balance=Some Balance 166 20 Creditworthy (0.8795181 0.1204819) *
- 3) Account.Balance=No Account 184 77 Creditworthy (0.5815217 0.4184783)
- 6) Duration.of.Credit.Month < 13 74 18 Creditworthy (0.7567568 0.2432432) *
- 7) Duration.of.Credit.Month >= 13 110 51 Non-Creditworthy (0.4636364 0.5363636)
- 14) Value.Savings.Stocks = < £100, £100-£1000 34 11 Creditworthy (0.6764706 0.3235294) *
- 15) Value.Savings.Stocks = None 76 28 Non-Creditworthy (0.3684211 0.6315789) *

Plots

- This is the report for our decision tree model.

Variable Importance



- We can see that the most important predictor variables are account balance, value savings stocks and duration of credit month.

Confusion Matrix					
	Creditworthy	Non-Creditworthy	Sum	Accuracy	
Actual	Creditworthy	225	28	253	89%
	Non-Creditworthy	49	48	97	49%
	Sum	274	76	350	78%
	Predicted				

- Confusion matrix for the decision tree.

ord Layout

Model Comparison Report						
Fit and error measures						
Model	Accuracy	F1	AUC	Accuracy_Creditworthy	Accuracy_Non-Creditworthy	
Decision_Tree	0.7467	0.8273	0.7054	0.8667	0.4667	
<p>Model: model names in the current comparison.</p> <p>Accuracy: overall accuracy, number of correct predictions of all classes divided by total sample number.</p> <p>Accuracy_[class name]: accuracy of Class [class name] is defined as the number of cases that are correctly predicted to be Class [class name] divided by the total number of cases that actually belong to Class [class name], this measure is also known as <i>recall</i>.</p> <p>AUC: area under the ROC curve, only available for two-class classification.</p> <p>F1: F1 score, $2 * \text{precision} * \text{recall} / (\text{precision} + \text{recall})$. The <i>precision</i> measure is the percentage of actual members of a class that were predicted to be in that class divided by the total number of cases predicted to be in that class. In situations where there are three or more classes, average precision and average recall values across classes are used to calculate the F1 score.</p>						
Confusion matrix of Decision_Tree						
		Actual_Creditworthy		Actual_Non-Creditworthy		
Predicted_Creditworthy		91		24		
Predicted_Non-Creditworthy		14		21		
Performance Diagnostic Plots						

- The overall accuracy for the decision tree is 74.67 % , which is a bit lower from the linear regression model. Also the model is biased for predicting non-creditworthy customers.

➔ **Forest Model :-**

Alteryx Designer x64 - project-template.yxmd - Browse (93)

9 records displayed, 2 fields, 85 KB

Table | Report | Profile

1 of 1 Fields | Records 1 to 9

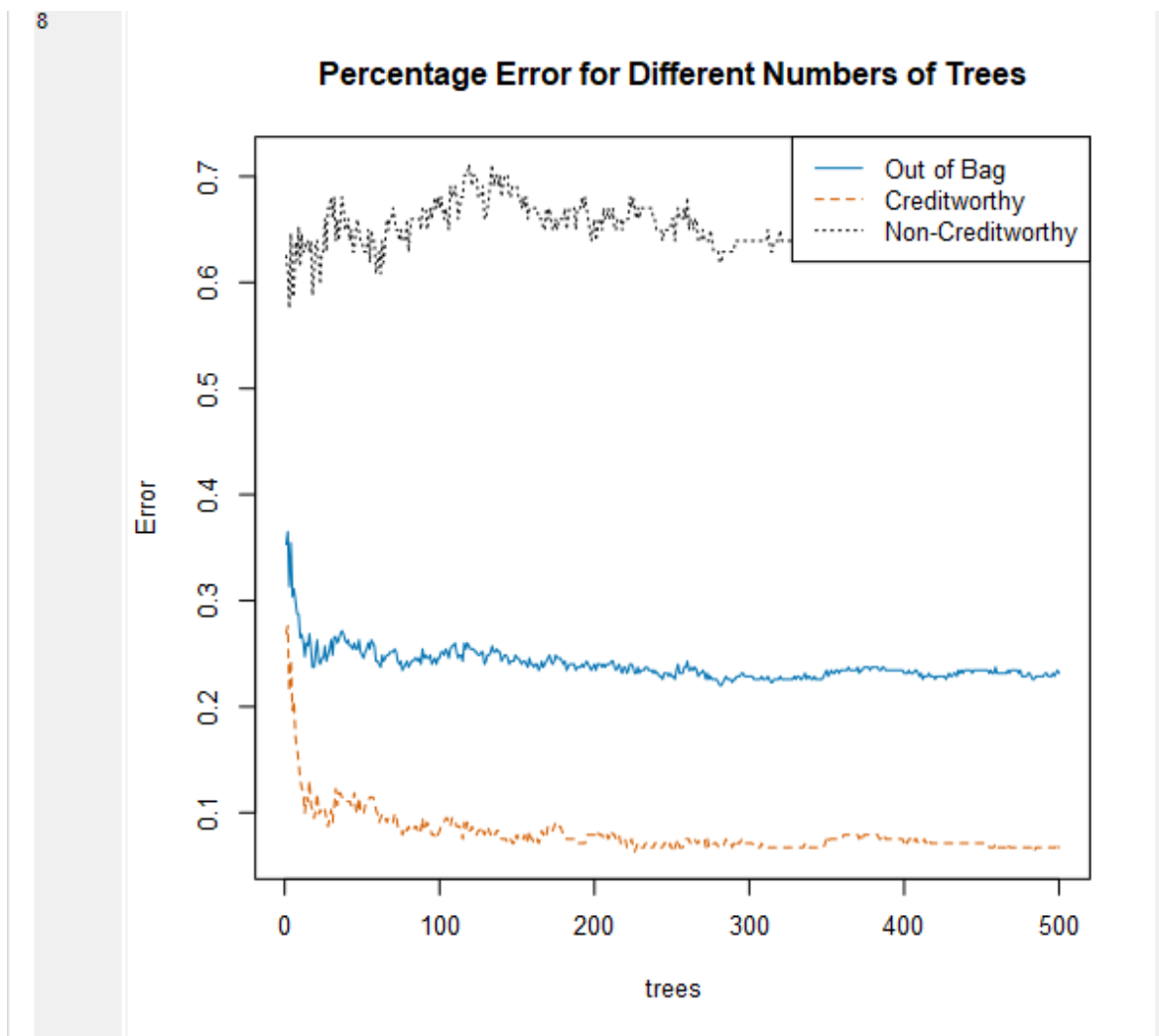
Record Report

- Basic Summary
- Call:
randomForest(formula = Credit.Application.Result ~ Account.Balance + Duration.of.Credit.Month + Payment.Status.of.Previous.Credit + Purpose + Credit.Amount + Value.Savings.Stocks + Length.of.current.employment + Instalment.per.cent + Most.valuable.available.asset + Age.years + Type.of.apartment + No.of.Credits.at.this.Bank, data = the.data, ntree = 500)
- Type of forest: classification
Number of trees: 500
Number of variables tried at each split: 3
- OOB estimate of the error rate: 36.3%
- Confusion Matrix:

	Classification Error	Creditworthy	Non-Creditworthy
Creditworthy	0.067	236	17
Non-Creditworthy	0.66	64	33

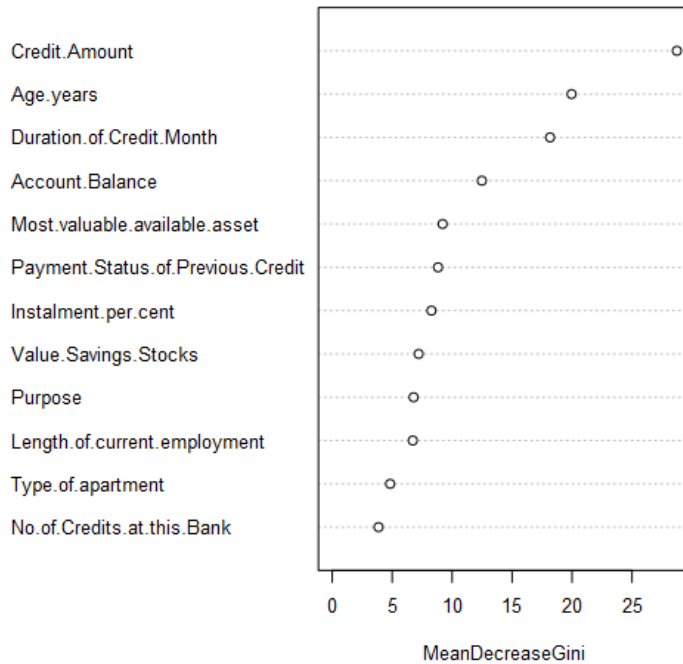
- Plots

- Report for the forest model



- Percent error for different types of trees.

Variable Importance Plot



- The most important predictor variables are credit amount, age-years and duration of credit month.

rd Layout

Model Comparison Report

Fit and error measures

Model	Accuracy	F1	AUC	Accuracy_Creditworthy	Accuracy_Non-Creditworthy
Forest_Model	0.7933	0.8681	0.7368	0.9714	0.3778

Model: model names in the current comparison.

Accuracy: overall accuracy, number of correct predictions of all classes divided by total sample number.

Accuracy_[class name]: accuracy of Class [class name] is defined as the number of cases that are **correctly** predicted to be Class [class name] divided by the total number of cases that actually belong to Class [class name], this measure is also known as *recall*.

AUC: area under the ROC curve, only available for two-class classification.

F1: F1 score, $2 * \text{precision} * \text{recall} / (\text{precision} + \text{recall})$. The *precision* measure is the percentage of actual members of a class that were predicted to be in that class divided by the total number of cases predicted to be in that class. In situations where there are three or more classes, average precision and average recall values across classes are used to calculate the F1 score.

Confusion matrix of Forest_Model

	Actual_Creditworthy	Actual_Non-Creditworthy
Predicted_Creditworthy	102	28
Predicted_Non-Creditworthy	3	17

Performance Diagnostic Plots

- We can see that the forest model offers 79.33% overall accuracy for the predictions. But this model is also not the best way to compute non-trustworthy customers as it offers 37% accuracy.

→ Boosted Model :-

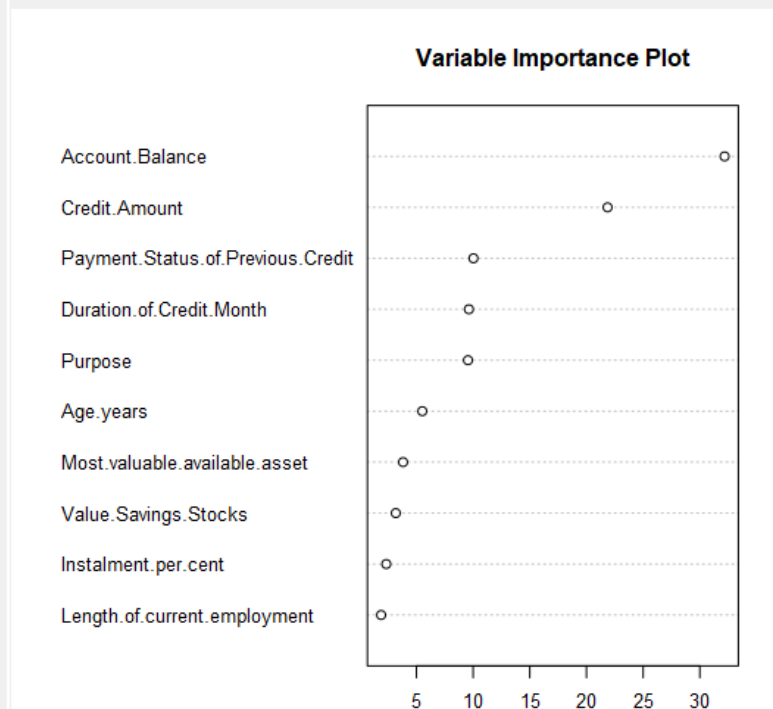
Loss function distribution: Bernoulli

Total number of trees used: 4000

Best number of trees based on 5-fold cross validation: 2036

2

Plots:



- We can see that the account balance, credit amount are the most significant ones for predicting the status of loan application.

Model Comparison Report					
Fit and error measures					
Model	Accuracy	F1	AUC	Accuracy_Creditworthy	Accuracy_Non-Creditworthy
Boosted_Model	0.7867	0.8632	0.7524	0.9619	0.3778
<p>Model: model names in the current comparison.</p> <p>Accuracy: overall accuracy, number of correct predictions of all classes divided by total sample number.</p> <p>Accuracy_[class name]: accuracy of Class [class name] is defined as the number of cases that are correctly predicted to be Class [class name] divided by the total number of cases that actually belong to Class [class name], this measure is also known as <i>recall</i>.</p> <p>AUC: area under the ROC curve, only available for two-class classification.</p> <p>F1: F1 score, $2 * \text{precision} * \text{recall} / (\text{precision} + \text{recall})$. The <i>precision</i> measure is the percentage of actual members of a class that were predicted to be in that class divided by the total number of cases predicted to be in that class. In situations where there are three or more classes, average precision and average recall values across classes are used to calculate the F1 score.</p>					
Confusion matrix of Boosted_Model					
	Actual_Creditworthy		Actual_Non-Creditworthy		
Predicted_Creditworthy	101		28		
Predicted_Non-Creditworthy	4		17		

- Boosted model offers 78.67 % overall accuracy for predicting the status of loan application, but it is also biased for predicting the non-creditworthy customers.

Step 4: Writeup

Decide on the best model and score your new customers. For reviewing consistency, if Score_Creditworthy is greater than Score_NonCreditworthy, the person should be labeled as "Creditworthy"

Write a brief report on how you came up with your classification model and write down how many of the new customers would qualify for a loan. (250 word limit)

Answer these questions:

- Which model did you choose to use? Please justify your decision using **all** of the following techniques. Please only use these techniques to justify your decision:
 ➔ By using the union and model comparison tools , I came up with the following report.

rd Layout

Model Comparison Report					
Fit and error measures					
Model	Accuracy	F1	AUC	Accuracy_Creditworthy	Accuracy_Non-Creditworthy
Decision_Tree	0.7467	0.8273	0.7054	0.8667	0.4667
Forest_Model	0.7933	0.8681	0.7368	0.9714	0.3778
Boosted_Model	0.7867	0.8632	0.7524	0.9619	0.3778
Stepwise	0.7600	0.8364	0.7306	0.8762	0.4889

Model: model names in the current comparison.

Accuracy: overall accuracy, number of correct predictions of all classes divided by total sample number.

Accuracy_[class name]: accuracy of Class [class name] is defined as the number of cases that are **correctly** predicted to be Class [class name] divided by the total number of cases that actually belong to Class [class name], this measure is also known as *recall*.

AUC: area under the ROC curve, only available for two-class classification.

F1: F1 score, $2 * \text{precision} * \text{recall} / (\text{precision} + \text{recall})$. The *precision* measure is the percentage of actual members of a class that were predicted to be in that class divided by the total number of cases predicted to be in that class. In situations where there are three or more classes, average precision and average recall values across classes are used to calculate the F1 score.

- ➔ We can see that the forest model is offering the highest accuracy as 79.33 %.
- ➔ The accuracy for predicting the creditworthy customers is 97.14 % ,although accuracy for predicting non-creditworthy customers is 37.78 % , but our main motive is to find customers which are trustworthy .

Records 1 to 8

Confusion matrix of Boosted_Model		
	Actual_Creditworthy	Actual_Non-Creditworthy
Predicted_Creditworthy	101	28
Predicted_Non-Creditworthy	4	17

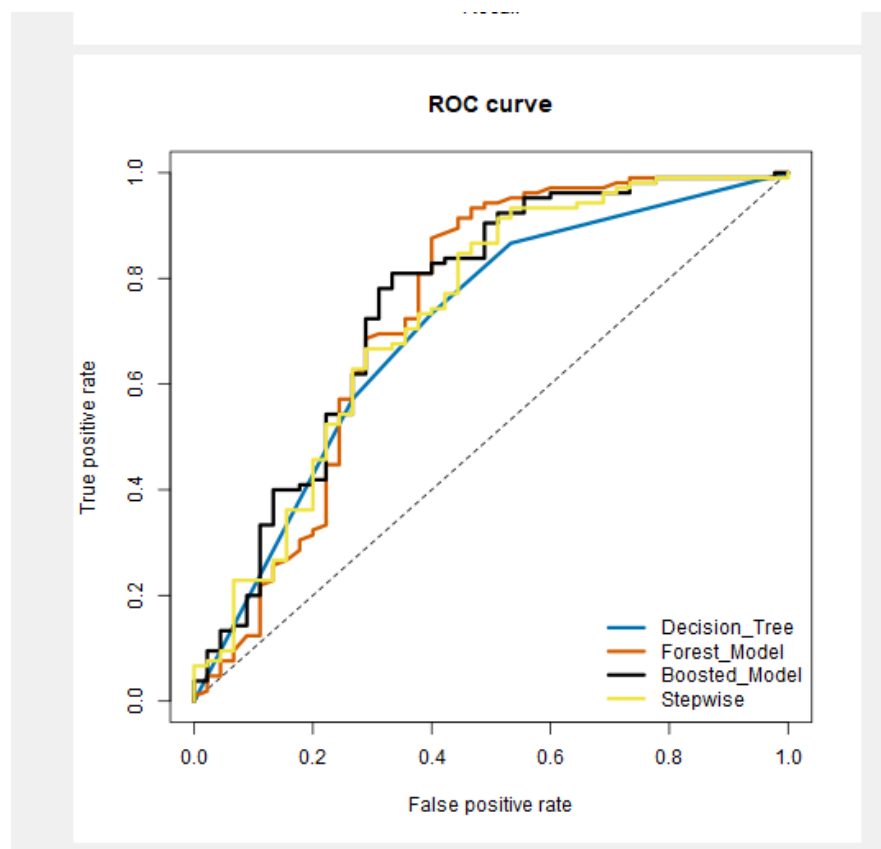
Confusion matrix of Decision_Tree		
	Actual_Creditworthy	Actual_Non-Creditworthy
Predicted_Creditworthy	91	24
Predicted_Non-Creditworthy	14	21

Confusion matrix of Forest_Model		
	Actual_Creditworthy	Actual_Non-Creditworthy
Predicted_Creditworthy	102	28
Predicted_Non-Creditworthy	3	17

Confusion matrix of Stepwise		
	Actual_Creditworthy	Actual_Non-Creditworthy
Predicted_Creditworthy	92	23
Predicted_Non-Creditworthy	13	22

ENGLISH (INDIA)

→ We can see the confusion matrices of all the models and can realize that forest model is the best among these four.



- We can see from the ROC curve that the forest model is performing the best.

- How many individuals are creditworthy?

Results - Browse (103) - Input

2 of 2 Fields | Cell Viewer | 1 record displayed, 1196 bytes

Record #1, Field Sum_Credit_App_Result_Creditworthy (Double)

Record #	Sum_Credit_App_Result_Creditworthy	Sum_Credit_App_Result_Non-Creditworthy
1	408	92

➔ After scoring the model and predicting data, I analysed that 408 customers are creditworthy and 92 customers are non-creditworthy.

The workflow for my work:-

