

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

1. What decisions needs to be made?

ANSWER: Out of a list of new 500 loan applicants, We need to decide , the number of people who are creditworthy.

2. What data is needed to inform those decisions?

ANSWER:We need Financial records of existing customers and of that of the possible clients. Besides , we need Income data, previous credit history , default history , net worth etc and other such inputs which impact a person's creditworthiness.

3. What kind of model (Continuous, Binary, Non-Binary, Time-Series) do we need to use to help make these decisions?

ANSWER:We can use a Binary models like a Logistic regression & non binary models like decision trees , boosted model & Forest model to help make these decisions.

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:

- *For numerical data fields, are there any fields that highly-correlate with each other? The correlation should be at least .70 to be considered "high".*
- *Are there any missing data for each of the data fields? Fields with a lot of missing data should be removed*
- *Are there only a few values in a subset of your data field? Does the data field look very uniform (there is only one value for the entire field?). This is called "low variability" and you should remove fields that have low variability. Refer to the "Tips" section to find examples of data fields with low-variability.*

- Your clean data set should have 13 columns where the Average of **Age Years** should be 36 (rounded up)

Note: For the sake of consistency in the data cleanup process, impute data using the average of the entire data field instead of removing a few data points. (100 word limit)

Note: For students using software other than Alteryx, please format each variable as:

Variable	Data Type
Credit-Application-Result	String
Account-Balance	String
Duration-of-Credit-Month	Double
Payment-Status-of-Previous-Credit	String
Purpose	String
Credit-Amount	Double
Value-Savings-Stocks	String
Length-of-current-employment	String
Instalment-per-cent	Double
Guarantors	String
Duration-in-Current-address	Double
Most-valuable-available-asset	Double
Age-years	Double
Concurrent-Credits	String
Type-of-apartment	Double
No-of-Credits-at-this-Bank	String
Occupation	Double
No-of-dependents	Double
Telephone	Double
Foreign-Worker	Double

To achieve consistent results reviewers expect.

Answer this question:

1. In your cleanup process, which fields did you remove or impute? Please justify why you removed or imputed these fields. Visualizations are encouraged.

ANSWER:

PREDICTOR VARIABLE REMOVED

Duration in Current Address

Foreign Worker

Concurrent credits

No. of dependents

REASON FOR REMOVAL

Too much missing data

Logical & Variable results are too skewed in one way.

Only one value in the variable

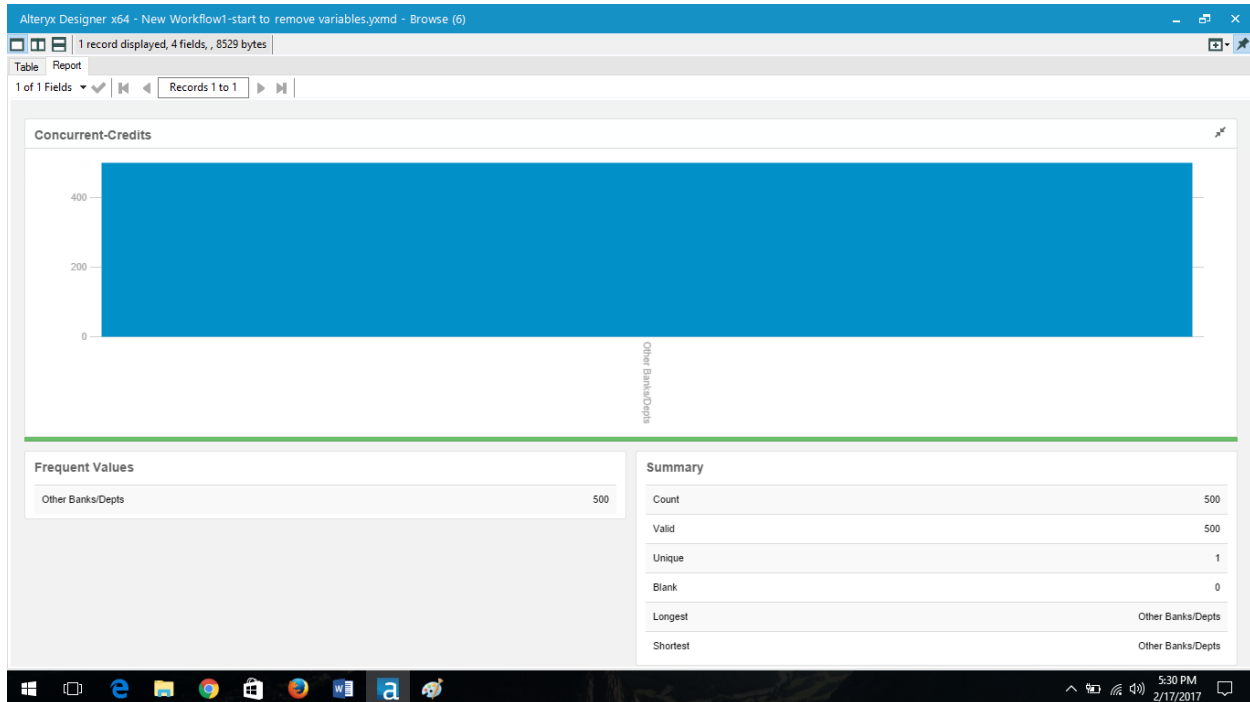
Logical. Has no impact on a person's creditworthiness.

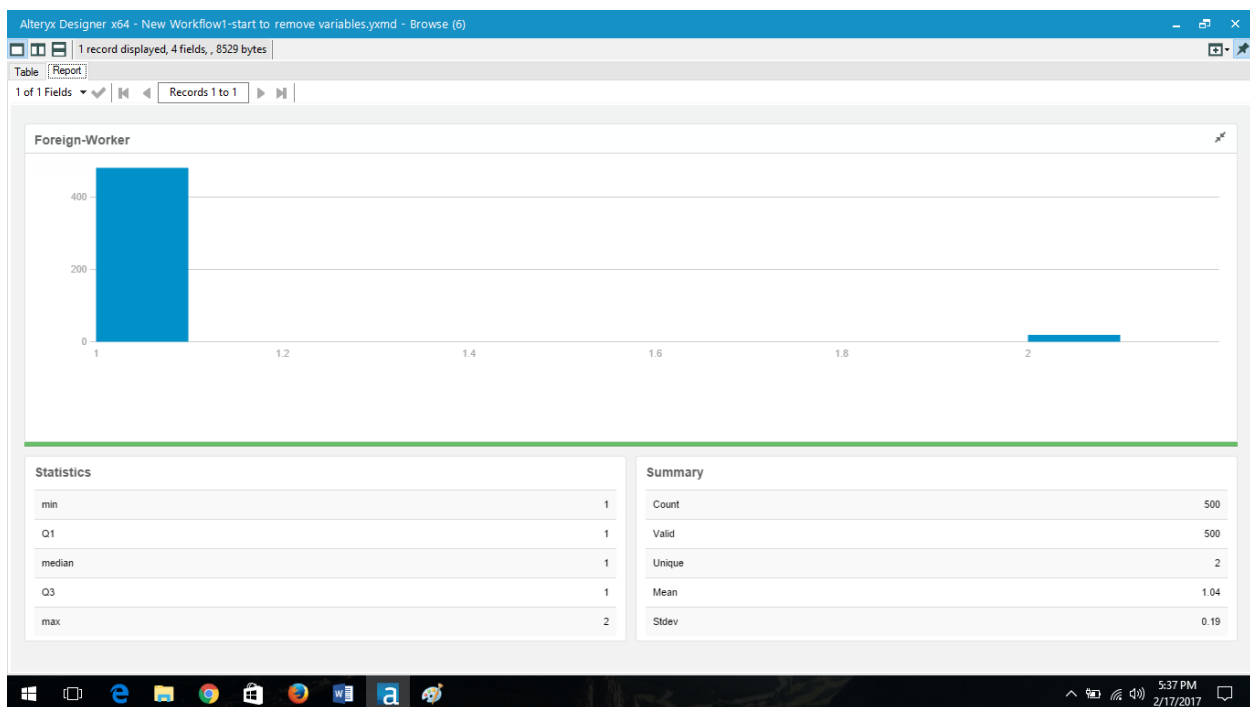
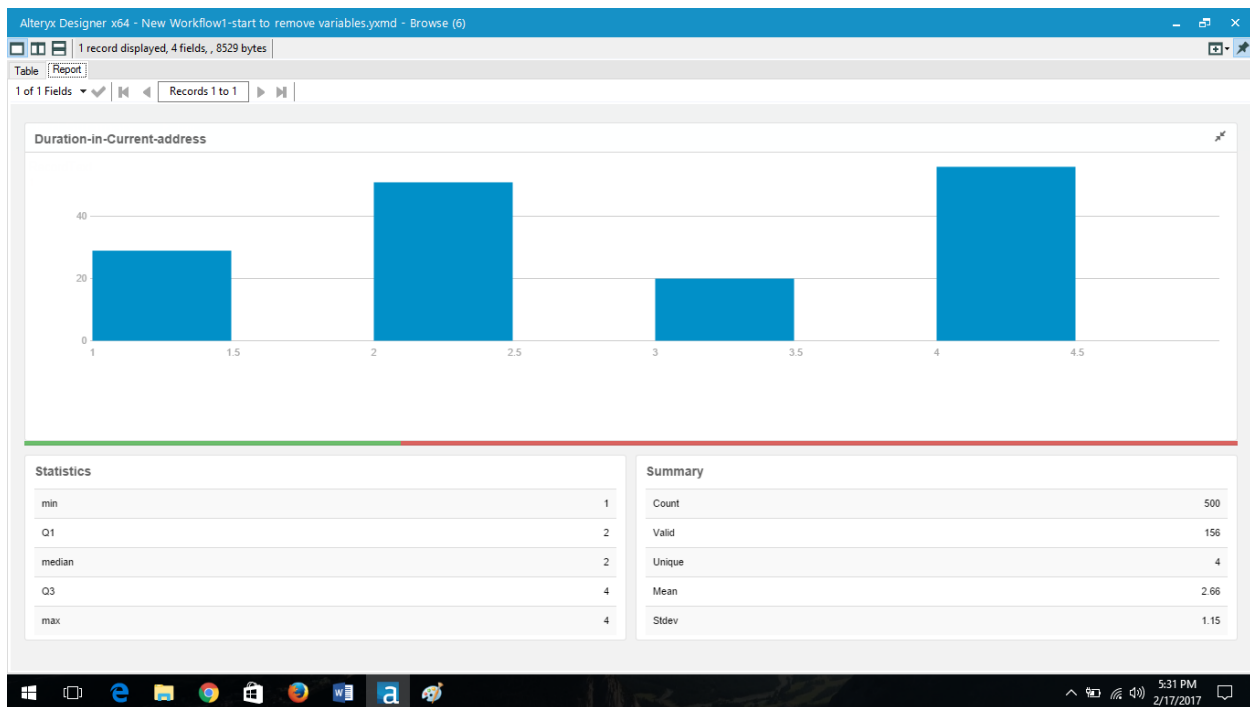
Guarantors
Telephone
Occupation

Variable results are too skewed in one way.
Logical. Has no impact on a person's creditworthiness
Only one value in the variable

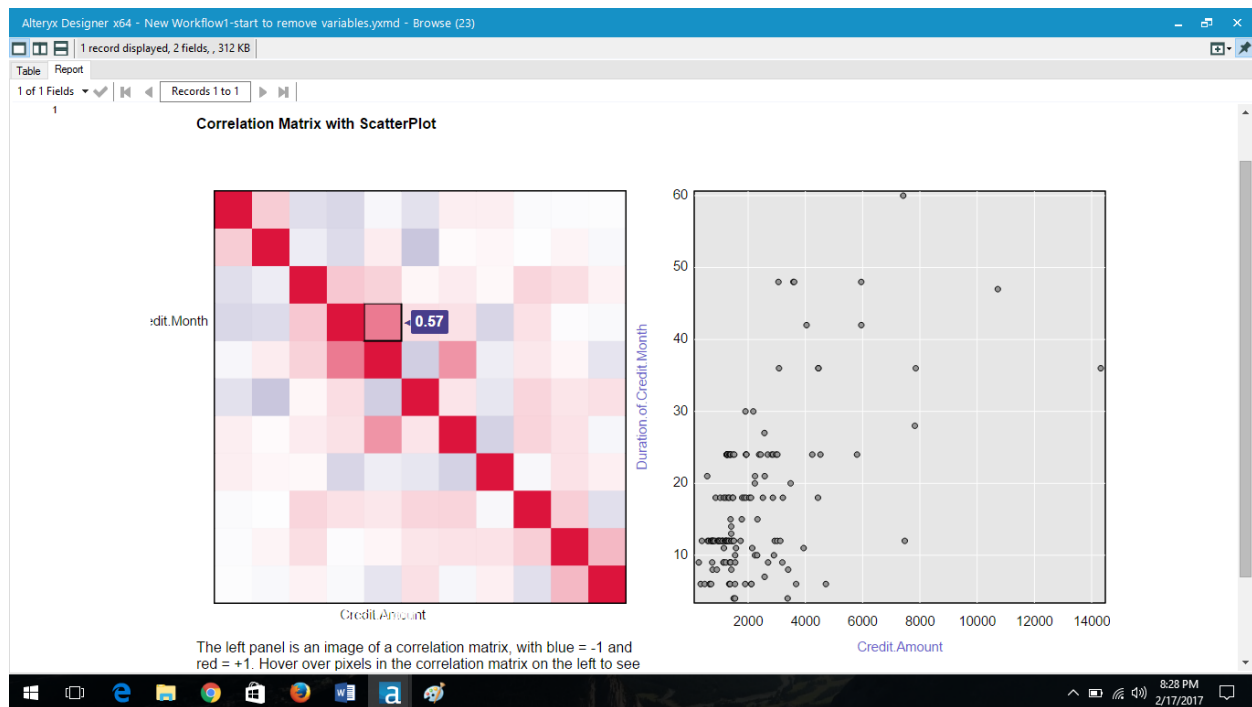
I also imputed the missing data in the "Age-years" field as only 2% data was missing and imputation would be the better option rather than removal, as in this case it would help build a model that would be closer to the reality. For this to happen, the missing Age fields would have to be filled in with the "average with null" values derived using the summarize tool and planted in the table using the formula tool.

Below are 3 sample visualizations that support the above answer:





I did not find numerical fields that highly correlate with each other(More than 0.70), Here is the scatterplot image below:



The maximum co-relation was between the fields “Credit amount” and Duration of credit month”. It is 0.57, as can be seen from the image above.

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:*

Confusion matrix of BoostedModel		
	Actual_Creditworthy	Actual_Non-Creditworthy
Predicted_Creditworthy	100	27
Predicted_Non-Creditworthy	5	18

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

Confusion matrix of ForestModel		
	Actual_Creditworthy	Actual_Non-Creditworthy
Predicted_Creditworthy	102	24
Predicted_Non-Creditworthy	3	21

Confusion matrix of Stepwise_Model		
	Actual_Creditworthy	Actual_Non-Creditworthy
Predicted_Creditworthy	92	32
Predicted_Non-Creditworthy	13	13

1. Which predictor variables are significant or the most important? Please show the p-values or variable importance charts for all of your predictor variables.
2. 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)

ANSWER:

MODEL 1: LOGISTIC REGRESSION & STEPWISE LOGISTIC REGRESSION

I first ran the logistic regression with all the possible predictor variables and I then ran the stepwise logistic regression tool so that it automatically selects the most significant variables. (Screenshots below)

Alteryx Designer x64 - New Workflow1-start to remove variables.yxmd - Browse (10)

13 records displayed, 2 fields, 98 KB

Table Report

1 of 1 Fields 1 of 1 Fields

Records 1 to 10

2 Basic Summary

3 Call:
glm(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, family = binomial(logit), data = the.data)

4 Deviance Residuals:

	Min	1Q	Median	3Q	Max
	-1.7160	-0.5439	-0.2777	-0.0297	3.0590

6 Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.0952041	1.6883664	-1.2410	0.21462
Account.BalanceSome Balance	-1.4981832	0.5596242	-2.6771	0.00743 **
Duration.of.Credit.Month	-0.0605157	0.0291709	-2.0745	0.03803 *
Payment.Status.of.Previous.CreditPaid Up	2.4391308	0.8251711	2.9559	0.00312 **
Payment.Status.of.Previous.CreditSome Problems	4.1596349	1.2241003	3.3981	0.00068 ***
PurposeNew car	0.1234228	0.8995404	0.1372	0.89087
PurposeOther	0.8729187	1.4503138	0.6019	0.54725
PurposeUsed car	-3.3643929	1.2298576	-2.7356	0.00623 **
Credit.Amount	0.0002066	0.0001198	1.7243	0.08465 .
Value.Savings.StocksNone	-1.0389796	0.9093136	-1.1426	0.25321
Value.Savings.Stocks£100-£1000	-2.1293301	1.0862684	-1.9602	0.04997 *
Length.of.current.employment4-7 yrs	1.9385821	1.0139289	1.9120	0.05588 .
Length.of.current.employment< 1yr	2.1868718	0.8446627	2.5890	0.00962 **
Instalment.per.cent	0.3978110	0.2450160	1.6236	0.10446
Most.valuable.available.asset	0.3122585	0.2912249	1.0722	0.28362
Age.years	0.0057132	0.0291564	0.1960	0.84465
Type.of.apartment	-1.2671259	0.5758109	-2.2006	0.02776 *
No.of.Credits.at.this.BankMore than 1	0.6878227	0.7424236	0.9265	0.35421

Alteryx Designer x64 - New Workflow1-start to remove variables.yxmd - Browse (20)

13 records displayed, 2 fields, 98 KB

Table Report

1 of 1 Fields

Records 1 to 10

Record Report

1 **Report for Logistic Regression Model Stepwise_Model**

2 **Basic Summary**

3 Call:
glm(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 + Type.of.apartment, family = binomial(logit), data = the.data)

4 Deviance Residuals:

5

	Min	1Q	Median	3Q	Max
	-1.7810	-0.5361	-0.2701	-0.0336	2.7220

6 Coefficients:

7

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-1.1156910	1.438398	-0.7756	0.43796
Account.BalanceSome Balance	-1.3443700	0.535371	-2.5111	0.01204 *
Duration.of.Credit.Month	-0.0552386	0.027919	-1.9785	0.04787 *
Payment.Status.of.Previous.CreditPaid Up	1.8795621	0.608212	3.0903	0.002 **
Payment.Status.of.Previous.CreditSome Problems	3.9043584	1.128594	3.4595	0.00054 ***
PurposeNew car	0.2792418	0.842235	0.3315	0.74023
PurposeOther	1.0536803	1.346944	0.7823	0.43405
PurposeUsed car	-3.2750986	1.195082	-2.7405	0.00613 **
Credit.Amount	0.0002177	0.000117	1.8610	0.06274 .
Value.Savings.StocksNone	-1.1162921	0.813206	-1.3727	0.16984
Value.Savings.Stocks£100-£1000	-2.1297611	0.999819	-2.1301	0.03316 *
Length.of.current.employment4-7 yrs	2.0574502	0.951306	2.1628	0.03056 *
Length.of.current.employment< 1yr	2.1112581	0.815857	2.5878	0.00966 **
Instalment.per.cent	0.4070390	0.236959	1.7178	0.08584 .
Type.of.apartment	-1.1157071	0.521423	-2.1397	0.03238 *

Significance codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Based on the model comparison report of the estimated dataset and the validation samples, the result is as follows:

Alteryx Designer x64 - New Workflow1-start to remove variables.yxmd - Browse (28)

5 records displayed, 2 fields, 89 KB

Table Report

1 of 1 Fields

Records 1 to 5

Record Layout

1 **Model Comparison Report**

2 **Fit and error measures**

Model	Accuracy	F1	AUC	Accuracy_Creditworthy	Accuracy_Non-Creditworthy
Stepwise_Model	0.6182	0.7386	0.5551	0.7063	0.3333

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], number of samples that are **correctly** predicted to be Class [class name] divided by number of samples predicted to be Class [class name]
AUC: area under the ROC curve, only available for two-class classification.
F1: F1 score, precision * recall / (precision + recall)

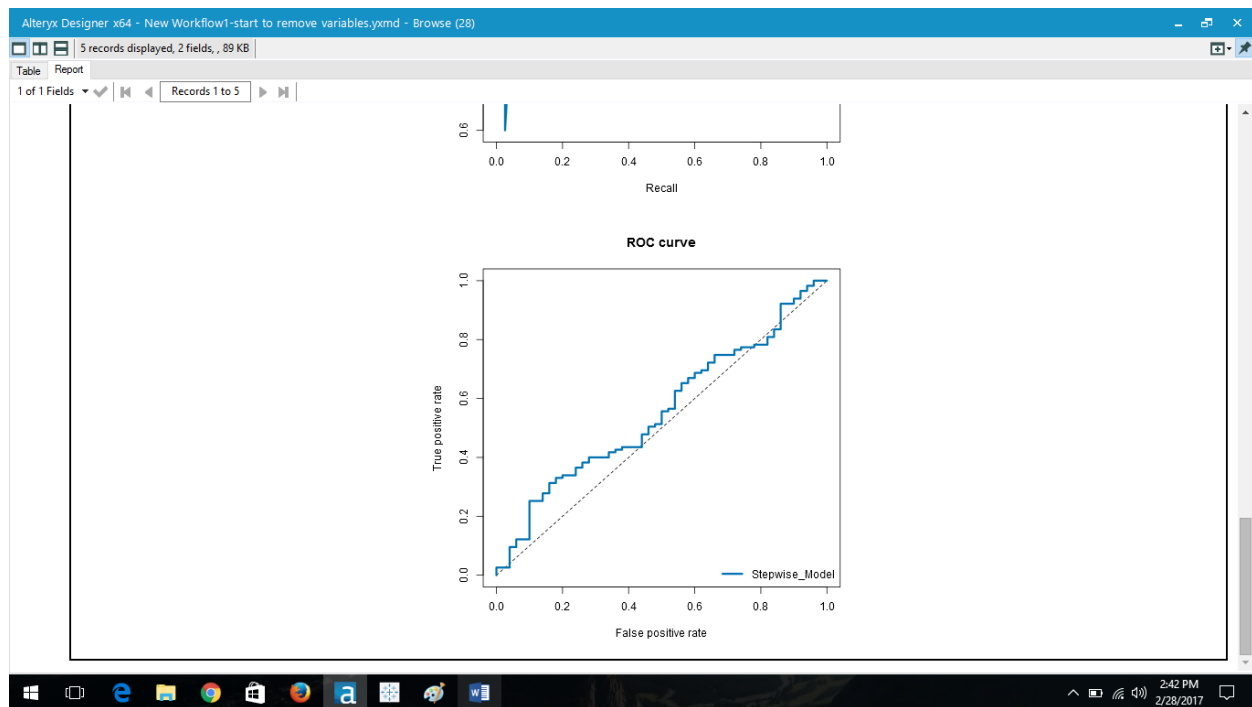
3 **Confusion matrix of Stepwise_Model**

	Actual_Creditworthy	Actual_Non-Creditworthy
Predicted_Creditworthy	89	37
Predicted_Non-Creditworthy	26	13

4 **Performance Diagnostic Plots**

5

Lift curve



To summarize: Accuracy of the model = 61.82%

Creditworthy predicted accurately: 70.63%

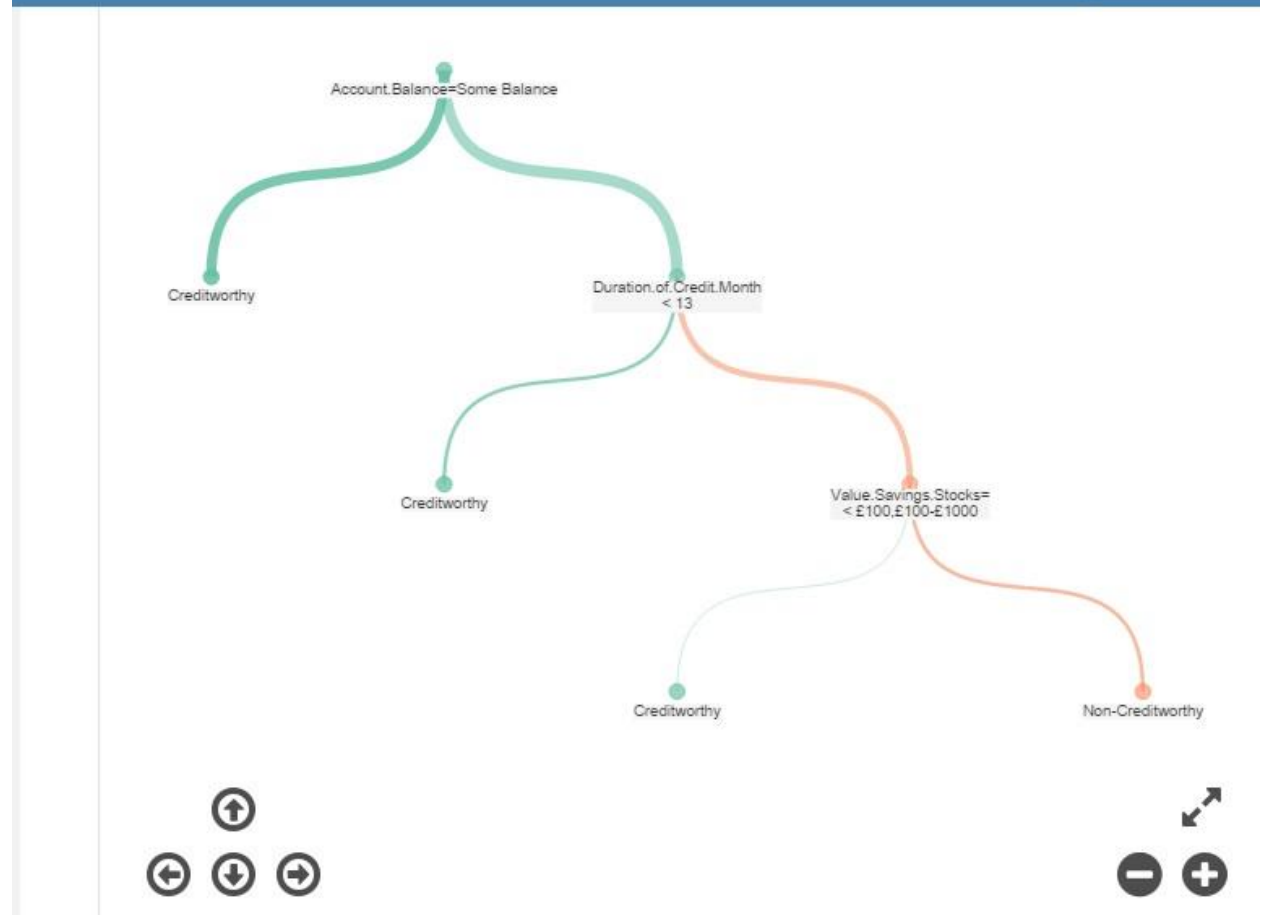
MODEL 2: DECISION TREE

Most Important variables as per this model = Account balance, Credit Amount, Duration of credit month, Value of Savings stocks amongst others as seen below

Variable Importance



Actual Decision tree



Confusion matrix

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

Alteryx Designer x64 - New Workflow1-start to remove variables.yxmd - Browse (34)																
5 records displayed, 2 fields, 75 KB																
Table Report																
1 of 1 Fields Records 1 to 5																
Record Layout																
1	Model Comparison Report															
2	<div>Fit and error measures</div> <table> <tr> <th>Model</th><th>Accuracy</th><th>F1</th><th>AUC</th><th>Accuracy_Creditworthy</th><th>Accuracy_Non-Creditworthy</th></tr> <tr> <td>DTree</td><td>0.7467</td><td>0.8273</td><td>0.7054</td><td>0.7913</td><td>0.6000</td></tr> </table> <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], number of samples that are correctly predicted to be Class [class name] divided by number of samples predicted to be Class [class name]</p> <p>AUC: area under the ROC curve, only available for two-class classification.</p> <p>F1: F1 score, precision * recall / (precision + recall)</p>				Model	Accuracy	F1	AUC	Accuracy_Creditworthy	Accuracy_Non-Creditworthy	DTree	0.7467	0.8273	0.7054	0.7913	0.6000
Model	Accuracy	F1	AUC	Accuracy_Creditworthy	Accuracy_Non-Creditworthy											
DTree	0.7467	0.8273	0.7054	0.7913	0.6000											
3	<div>Confusion matrix of DTree</div> <table> <tr> <th></th><th>Actual_Creditworthy</th><th>Actual_Non-Creditworthy</th></tr> <tr> <th>Predicted_Creditworthy</th><td>91</td><td>24</td></tr> <tr> <th>Predicted_Non-Creditworthy</th><td>14</td><td>21</td></tr> </table>					Actual_Creditworthy	Actual_Non-Creditworthy	Predicted_Creditworthy	91	24	Predicted_Non-Creditworthy	14	21			
	Actual_Creditworthy	Actual_Non-Creditworthy														
Predicted_Creditworthy	91	24														
Predicted_Non-Creditworthy	14	21														
4	Performance Diagnostic Plots															
5	<div>Lift curve</div>															

To summarize:

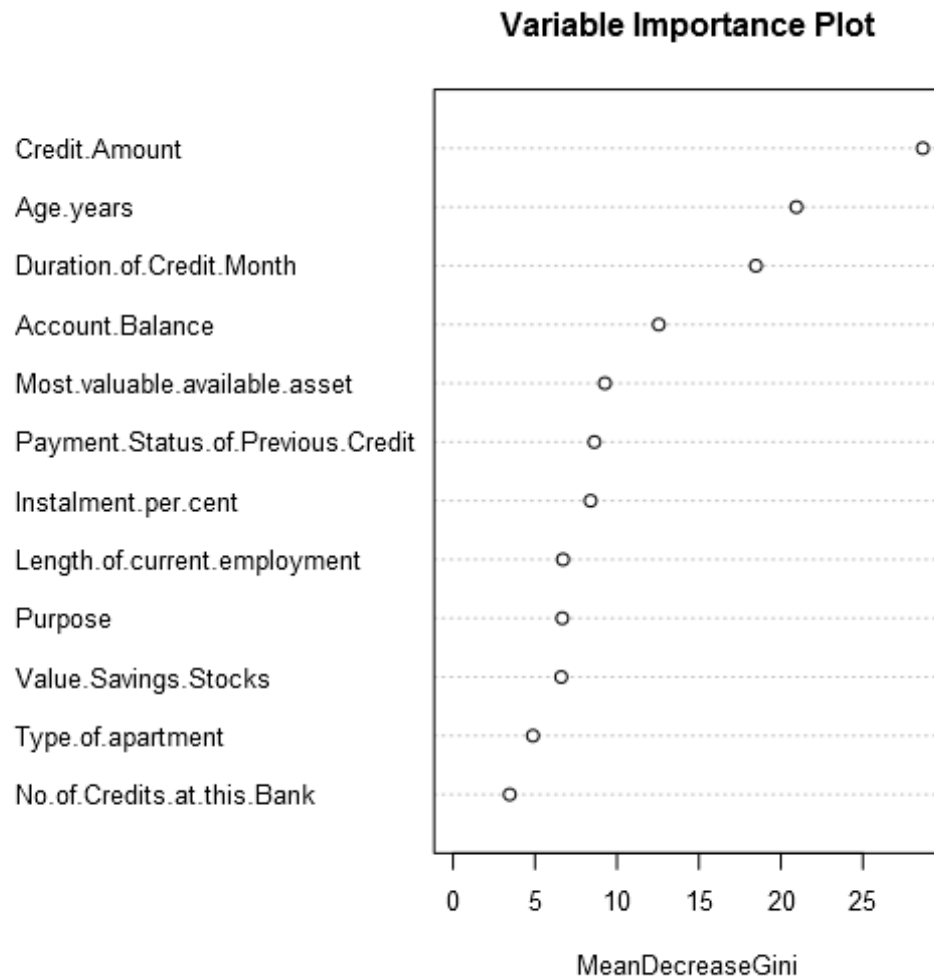
Accuracy of the model = 74.67%

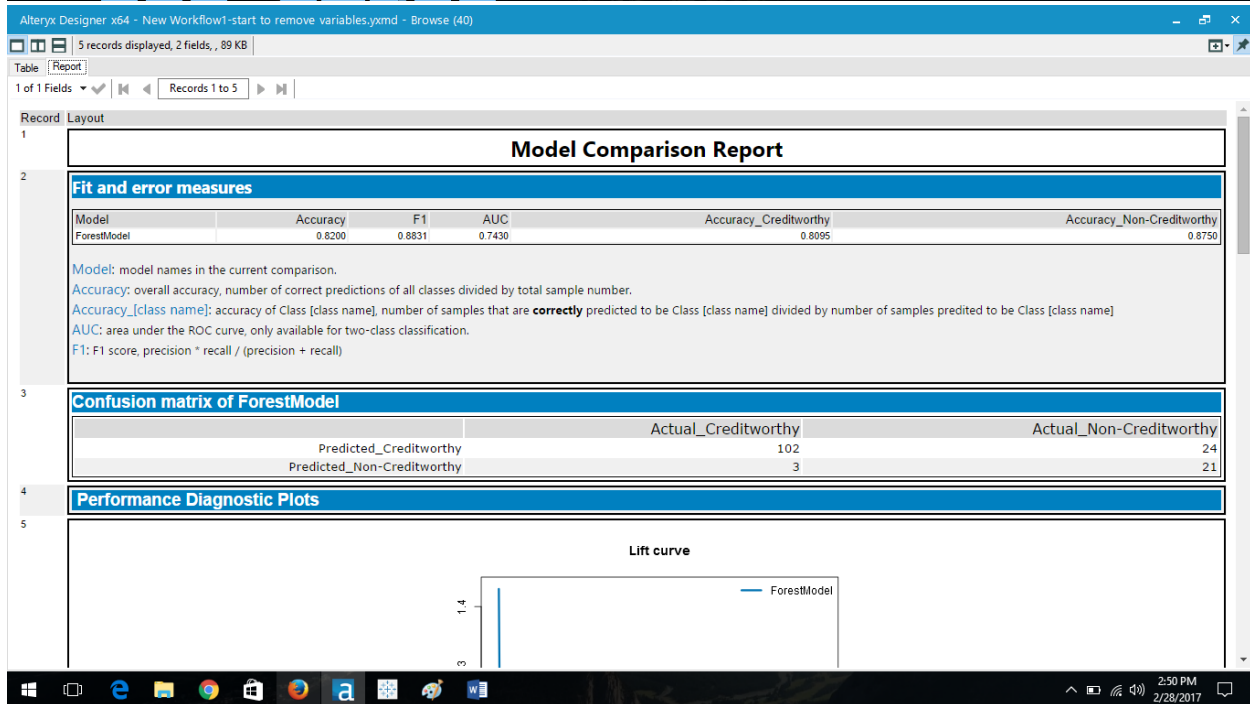
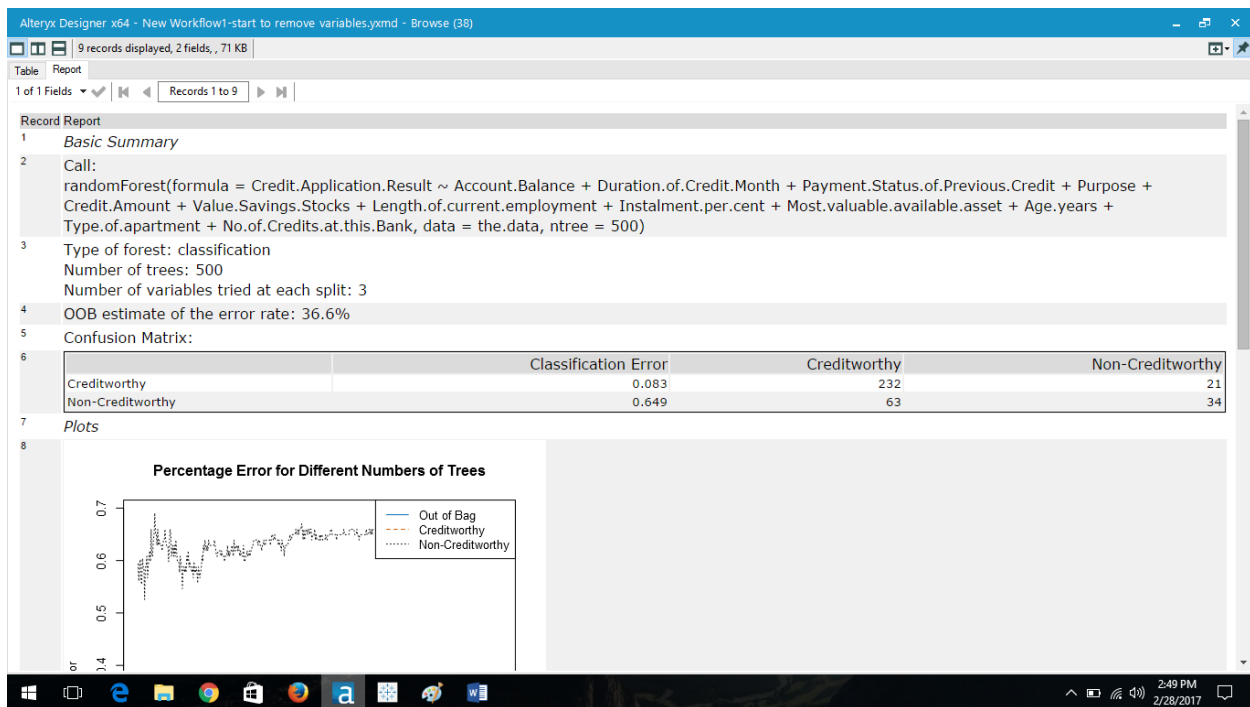
Creditworthy predicted accurately: 79.13%

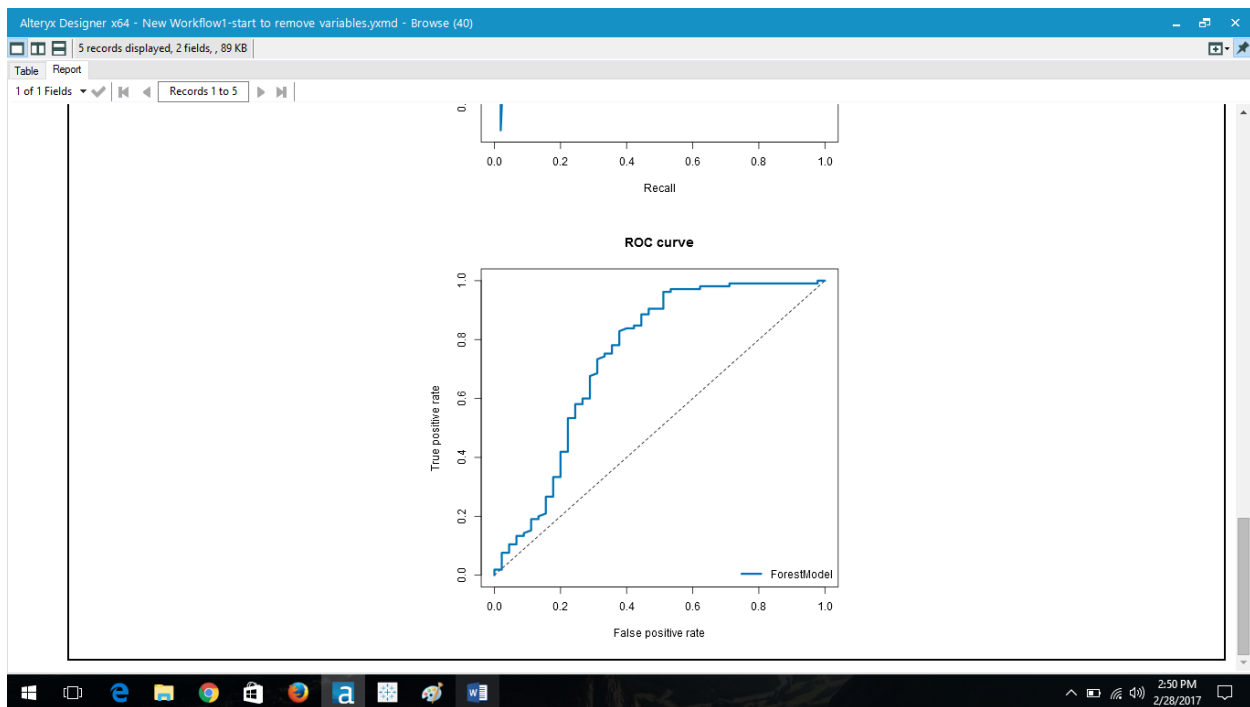
MODEL 3: Forest Model

Variables importance plot

Most Important variables – Credit Amount, Age years, Duration of credit month ,Account balance ,most valuable asset amongst others





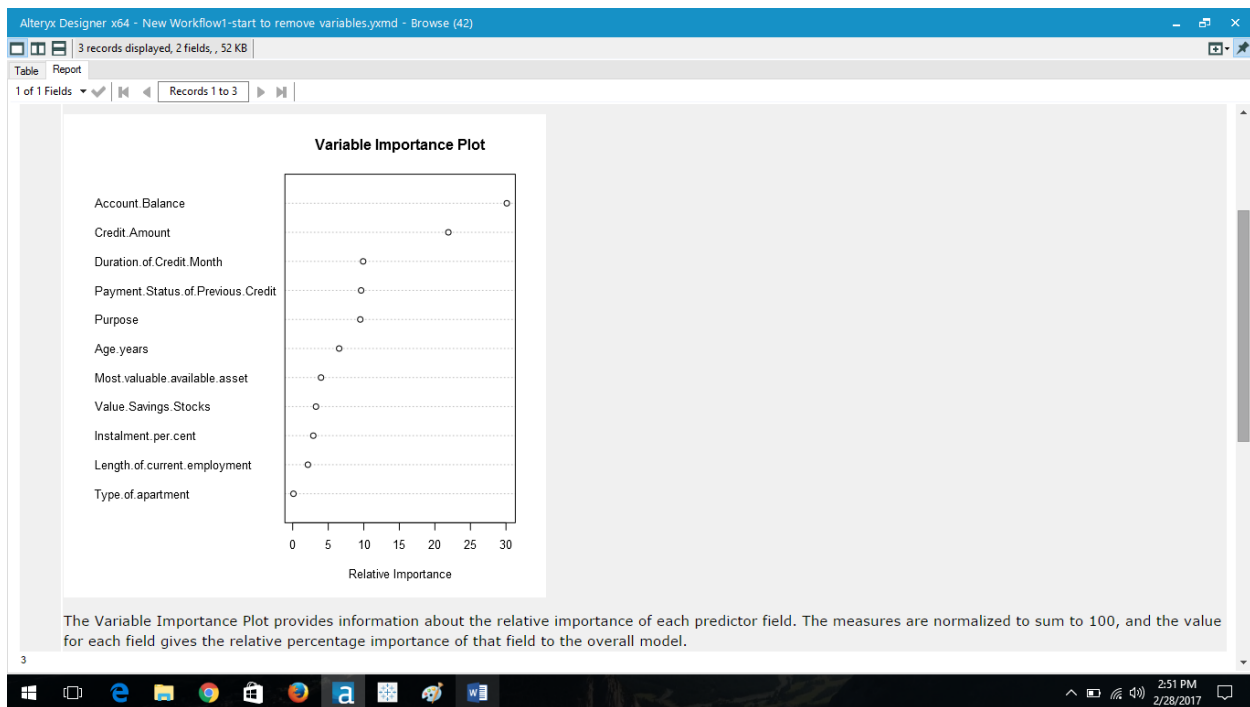


*To summarize: Accuracy of the model: 82%
Creditworthy predicted accurately: 80.95%*

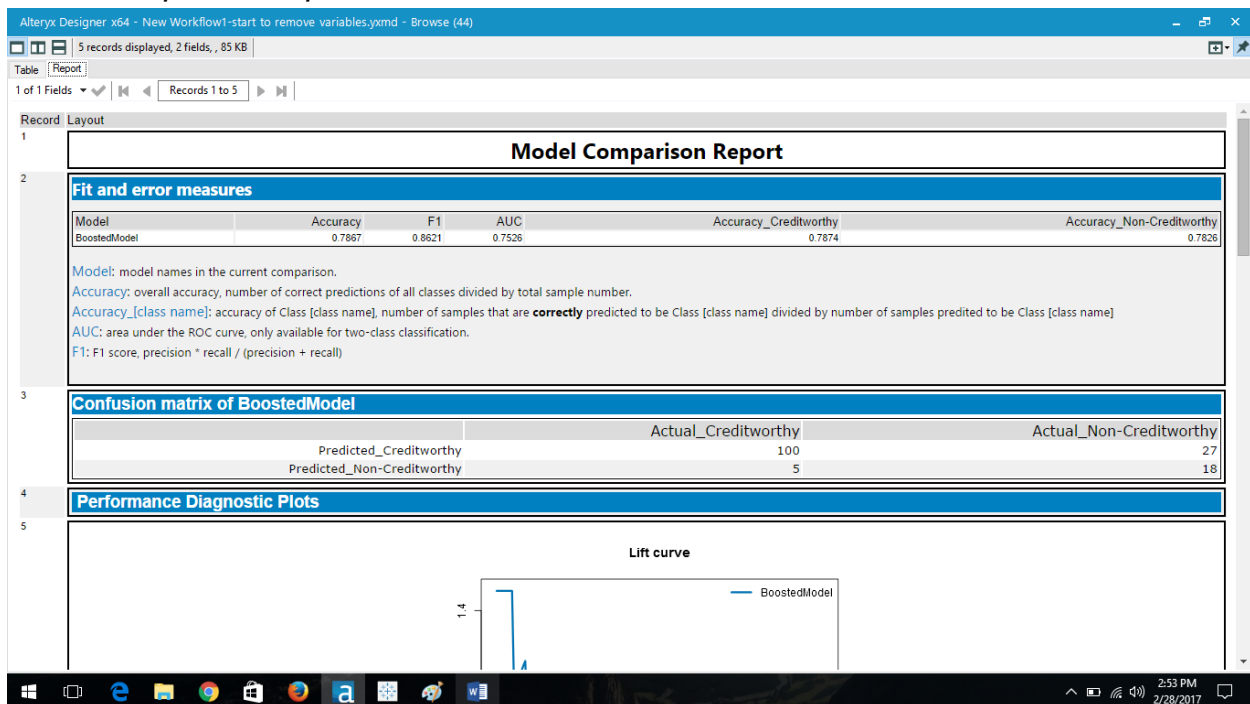
MODEL 4: Boosted Model

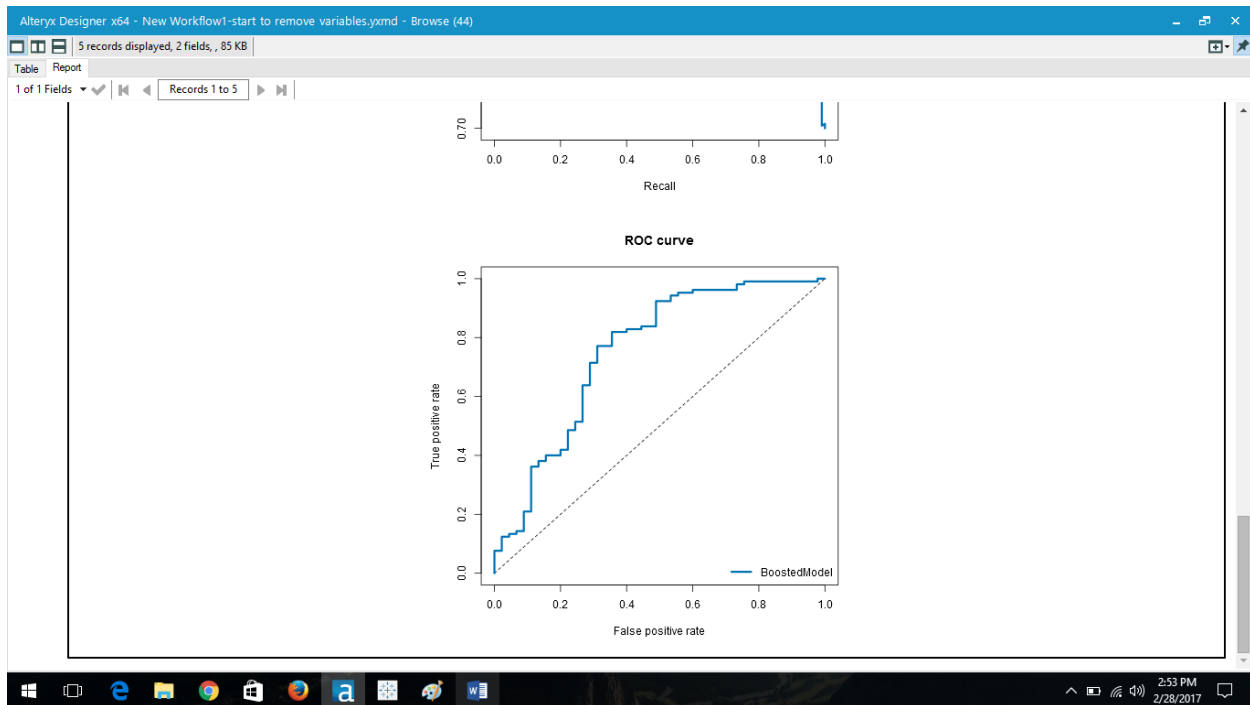
Variable importance plot

*Most Important Variables: Account balance, credit amount, Duration of credit month,
Payment status of Previous credit, Purpose amongst others*



Model comparison report





To summarize: Accuracy of the model: 78.67%

Creditworthy predicted accurately: 78.74%

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:

1. Which model did you choose to use? Please justify your decision using only the following techniques:
 - a. Overall Accuracy against your Validation set
 - b. Accuracies within "Creditworthy" and "Non-Creditworthy" segments
 - c. ROC graph
 - d. Bias in the Confusion Matrices

ANSWER:

Let's begin with the reports,

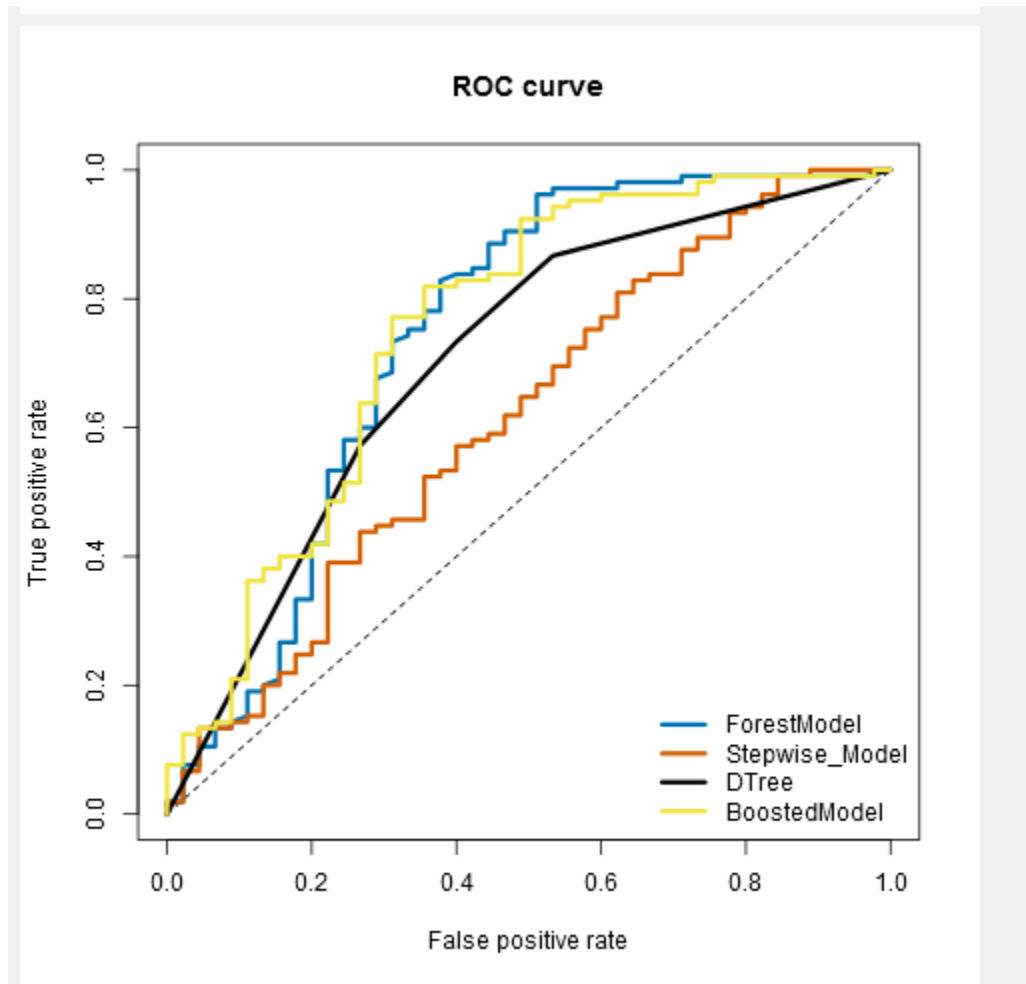
Model Comparison Report						
Fit and error measures						
Model	Accuracy	F1	AUC	Accuracy_Creditworthy	Accuracy_Non-Creditworthy	
ForestModel	0.8200	0.8831	0.7430	0.8095	0.8750	
Stepwise_Model	0.7000	0.8035	0.6142	0.7419	0.5000	
DTree	0.7467	0.8273	0.7054	0.7913	0.6000	
BoostedModel	0.7867	0.8621	0.7526	0.7874	0.7826	

Confusion matrix of BoostedModel			
	Actual_Creditworthy	Actual_Non-Creditworthy	
Predicted_Creditworthy	100	27	
Predicted_Non-Creditworthy	5	18	

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

Confusion matrix of ForestModel			
	Actual_Creditworthy	Actual_Non-Creditworthy	
Predicted_Creditworthy	102	24	
Predicted_Non-Creditworthy	3	21	

Confusion matrix of Stepwise_Model			
	Actual_Creditworthy	Actual_Non-Creditworthy	
Predicted_Creditworthy	92	32	
Predicted_Non-Creditworthy	13	13	



Based on the above reports, Considering the overall accuracy of the model & the accuracy of the “Creditworthy” applicants validated, I conclude that the Forest Model is the best fit in this case. Though the ROC curve is marginally better of the given boosted model (derived visually) and also through the AUC values on the above report, the forest model AUC values & ROC curves are very close, considering all the 3 factors together- Overall Accuracy, Accuracy of Creditworthy predictions derived through the confusion matrices and the ROC curve, I believe the Forest Model best fits this case.

To summarize the Forest Model

Overall accuracy – 82.00%

Creditworthy predicted accurately: 80.95%

Note: Remember that your boss only cares about prediction accuracy for Creditworthy and Non-Creditworthy segments.

2. How many individuals are creditworthy?

When we score the above model, 406 applicants out of the 500 are determined to be creditworthy

Before you Submit

Please check your answers against the requirements of the project dictated by the [rubric](#) here. Reviewers will use this rubric to grade your project.