# Project: Creditworthiness
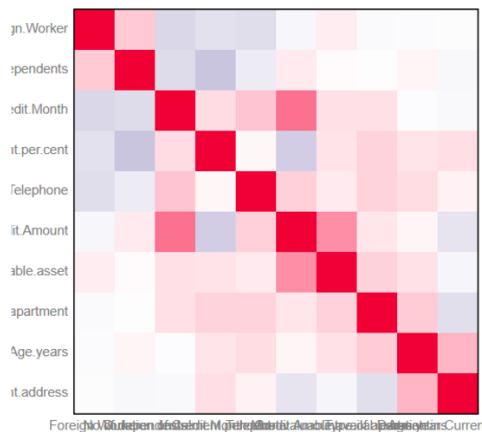
## Step 1: Business and Data Understanding

1. What decisions needs to be made?
   a. The objective is to determine whether new customers are creditworthy to the loan they applied.

2. What data is needed to inform those decisions?
   a. Data on past applications such as "Account-Balance" and "Credit-Amount" are needed to inform the decisions.

3. What kind of model (Continuous, Binary, Non-Binary, Time-Series) do we need to use to help make these decisions?
   a. Binary classification models like "Logistic Regression", "Decision Tree", "Forest Model", and "Boosted Model" are required to make these decisions.
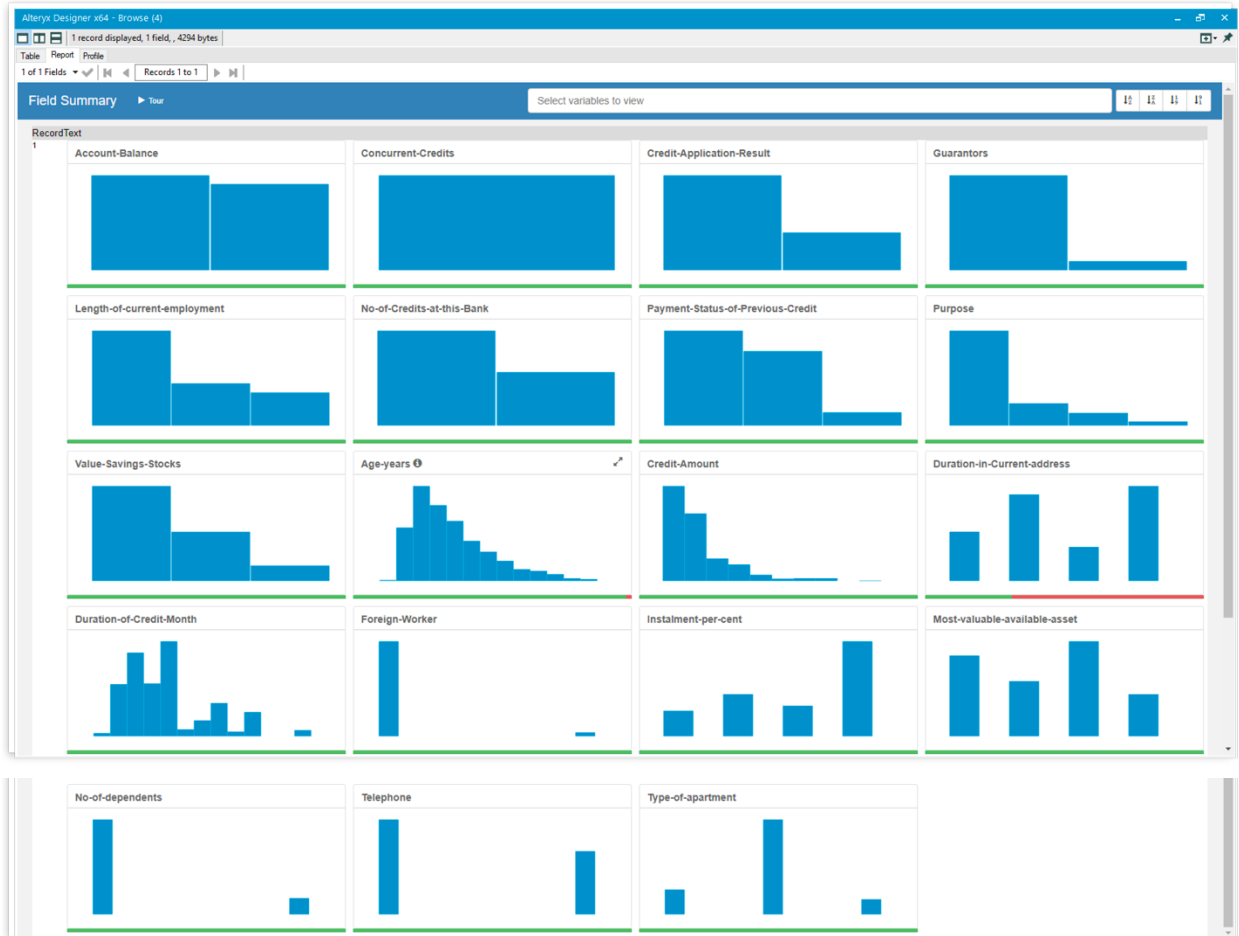
## Step 2: Building the Training Set

1. For the numerical data fields, an association analysis is performed to see the correlation. There are no numerical variables that are highly correlated with each other (correlation greater than 0.7).



Correlation Matrix with ScatterPlot

2. In your cleanup process, which fields did you remove or impute? Please justify why you removed or imputed these fields. Visualizations are encouraged.
   a. A field summary tool was performed on all the variables to identify any missing values. Out of all the variables, "Duration-in-Current-address" had 69% missing data and "Age-Years" had 2% missing data.
   b. "Duration-in-Current-address" should be removed since it has the most missing data.

c. Since "Age-Years" has few missing data, it should be imputed with median age.
d. In terms of "Low Variability", "Occupation" only had one value, "Concurrent-Credits", "Guarantors", "Foreign-Worker", "No-of-dependents", and "Telephone" has most of their data on one side indicating low variability. So, they should be removed.

# Step 3: Train your Classification Models

## 1. Logistic Regression:

a. **Stepwise:** The target variable used is "Credit-Application-Result" and selecting all the predictor variables expect the Credit-Application-Result. This resulted in the most significant variables with p-value less than 0.5 are: Account-Balance, Purpose, and Credit-Amount.



b. **Model Comparison:** The accuracy of the Stepwise model is 76%. Creditworthy is at 87%, higher than Non-Creditworthy which is at 48%.

## 2. Decision Tree:

a. The "Root Node Error" is at 2.7%, which indicates that 2.7% were predicted incorrectly.



b. There are three variables that have significant importance: Account-Balance, Value-Saving-Stocks, and Duration-of-Credit-Month.

c. In the Confusion Matrix, the Creditworthy accuracy is at 89%, while the Non-Creditworthy accuracy is at 49%. However, the overall accuracy is at 78%.

d. **Model Comparison:** The accuracy of the Decision Tree in the Model Comparison tool is 76%. Creditworthy accuracy is at 86% and the Non-Creditworthy accuracy is at 46%.



**Alteryx Designer x64 - Browse (19)**

Table | Report | Profile
1 of 1 Fields ▾ ✓ ◄ ◄  Records 1 to 5  ► ►

Record Layout

### Model Comparison Report

**Fit and error measures**

| Model | Accuracy | F1 | AUC | Accuracy_Creditworthy | Accuracy_Non-Creditworthy |
|---|---|---|---|---|---|
| DT_Credit | 0.7467 | 0.8273 | 0.7054 | 0.8667 | 0.4667 |

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 * precision * recall / (precision + 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 DT_Credit**

| | Actual_Creditworthy | Actual_Non-Creditworthy |
|---|---|---|
| Predicted_Creditworthy | 91 | 24 |
| Predicted_Non-Creditworthy | 14 | 21 |

*3. Forest Model:*

    a. The type of forest used is classification. The number of trees used, 500 (default). Number of splits, 3.

    b. The Out of Body (OOB) estimation error rate is at 35.8%, which is quite high.

    c. Compared to the OOB, the confusion matrix shows a better estimation error rate. Creditworthy at 8.7% and Non-Creditworthy at 62.9%



**d. Percentage error for different number of tress:**



    e. **Variable Importance:** The most important variables are: Credit-Amount, Age-Years, and Duration-of-Credit-Month which have plots on the top right side of the graph.

f. **Model Comparison:** The accuracy for the Forest model is at 80%. The Creditworthy accuracy at 96% and Non-Creditworthy at 44%.

**Model Comparison Report**

**Fit and error measures**

| Model | Accuracy | F1 | AUC | Accuracy_Creditworthy | Accuracy_Non-Creditworthy |
|---|---|---|---|---|---|
| F_Credit | 0.8067 | 0.8745 | 0.7365 | 0.9619 | 0.4444 |

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 * precision * recall / (precision + 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.
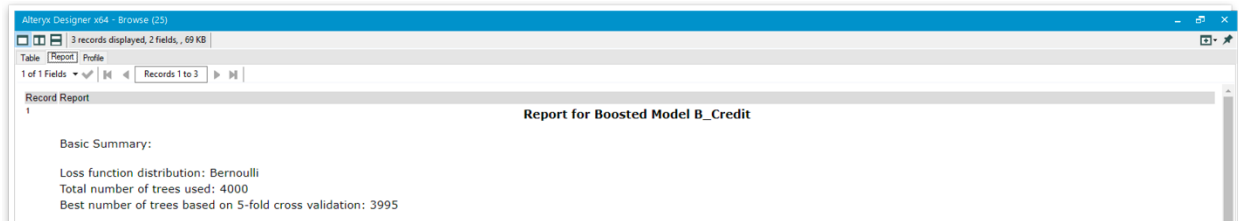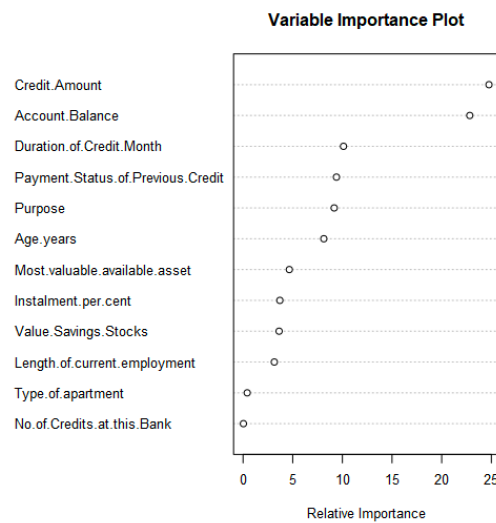
**Confusion matrix of F_Credit**

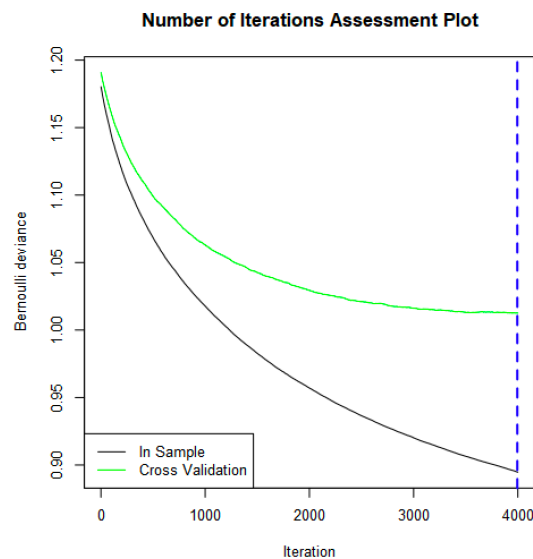| | Actual_Creditworthy | Actual_Non-Creditworthy |
|---|---|---|
| Predicted_Creditworthy | 101 | 25 |
| Predicted_Non-Creditworthy | 4 | 20 |

a. Since we are trying to predict whether to give the customers loan or not, there are only two outcomes, hence, Bernoulli.

Alteryx Designer x64 - Browse (25)

3 records displayed, 2 fields, , 69 KB

Table | Report | Profile

1 of 1 Fields | Records 1 to 3

Record Report
1

**Report for Boosted Model B_Credit**

Basic Summary:

Loss function distribution: Bernoulli
Total number of trees used: 4000
Best number of trees based on 5-fold cross validation: 3995

b. **Variable Importance:** The two most variables that are significant are: Credit-Amount and Account-Balance.

**Variable Importance Plot**

Credit.Amount
Account.Balance
Duration.of.Credit.Month
Payment.Status.of.Previous.Credit
Purpose
Age.years
Most.valuable.available.asset
Instalment.per.cent
Value.Savings.Stocks
Length.of.current.employment
Type.of.apartment
No.of.Credits.at.this.Bank

0   5   10   15   20   25

Relative Importance

c. **Number of Iteration Assessment:**

**Number of Iterations Assessment Plot**

Bernoulli deviance

— In Sample
— Cross Validation

0   1000   2000   3000   4000

Iteration

d. **Model Comparison:** The accuracy for the Boosted model is at 78%. Creditworthy accuracy is at 95% and Non-Creditworthy accuracy at 40%. The accuracies between Creditworthy and Non-Creditworthy is most biased.

5 records displayed, 2 fields, , 87 KB

Table  Report  Profile

1 of 1 Fields  ▾ ✔  |◀  ◀   Records 1 to 5   ▶  ▶|

Record Layout

**Model Comparison Report**

**Fit and error measures**

| Model | Accuracy | F1 | AUC | Accuracy_Creditworthy | Accuracy_Non-Creditworthy |
|---|---|---|---|---|---|
| B_Credit | 0.7867 | 0.8621 | 0.7526 | 0.9524 | 0.4000 |

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 * precision * recall / (precision + 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.
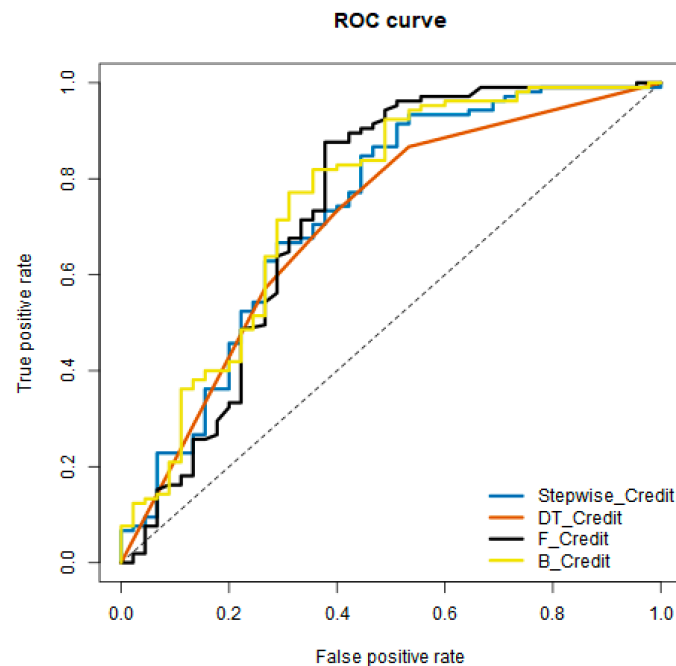
**Confusion matrix of B_Credit**

|  | Actual_Creditworthy | Actual_Non-Creditworthy |
|---|---|---|
| Predicted_Creditworthy | 100 | 27 |
| Predicted_Non-Creditworthy | 5 | 18 |

# Step 4: Writeup

## 1. *Model Comparison:*

    a. Out of all the four models, Forest Model has the highest accuracy at 80%. It also has highest Creditworthy accuracy at 96% and the least Non-Creditworthy accuracy at 44%.



    b. **ROC Curve:** The Forest Model was the first to reach the top in the ROC curve.



    c. The accuracies between Creditworthy and Non-Creditworthy lead to the least bias in the Confusion Matrix.

d. There is a total of 416 creditworthy customers.



Results - Browse (4) - Input

2 of 2 Fields ▼ ✓ | Cell Viewer ▼ | ↑ ↓ | 1 record displayed, 1098 bytes

| Record # | Sum_Score_Creditworthy | Sum_Score_Non-Creditworthy |
|---|---|---|
| 1 | 416 | 84 |