# Project: Creditworthiness

## Step 1: Business and Data Understanding

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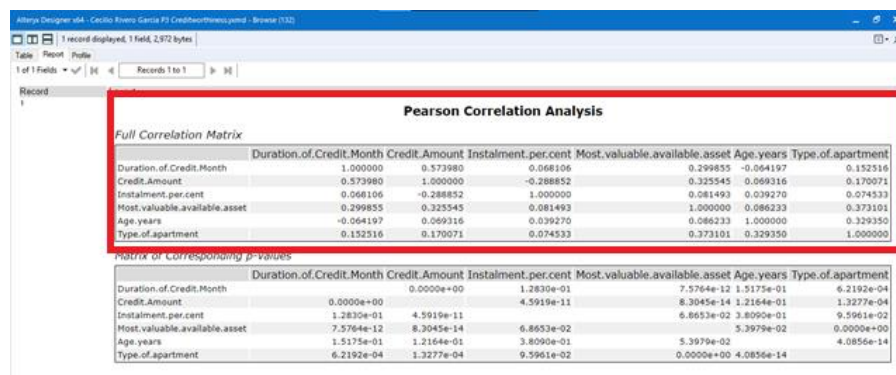
## Key Decisions:

Answer these questions

- What decisions needs to be made?

  It is required to find out if a loan application could be approved or rejected considering socio-economic data; this involves choosing the best predictive model (the one that shows the highest accuracy and least possible biasing).

- What data is needed to inform those decisions?

  The data to consider should be socio-economics data like: credit amount, account balance, duration of the credit, requester age and profession, etc. of previous loan applicants and their behavior: creditworthy or non-creditworthy.

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

  Because of the kind of answer involved (requester: creditworthy or non-creditworthy), the predictive model to use must be a Binary Model, like: Logistic Regression, Logistic Regression-Stepwise, Decision Tree, Forest Model or Boosted 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:*

- 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".

  Unfortunately, there is not any numerical data filed that show a highly-correlation with each other according to the Pearson Correlation Analysis.

- Are there any missing data for each of the data fields? Fields with a lot of missing data should be removed

  Two fields show up some missing data:
  - Duration-in-current-address with 69% of missing data. Because of the high amount of missing data, this field is not considered in the predictive model
  - Age-year with 2% of missing data. Because of the low amount of missing data, the records of this filed with null value were filled with the amount 33 (the median of Age-year field)
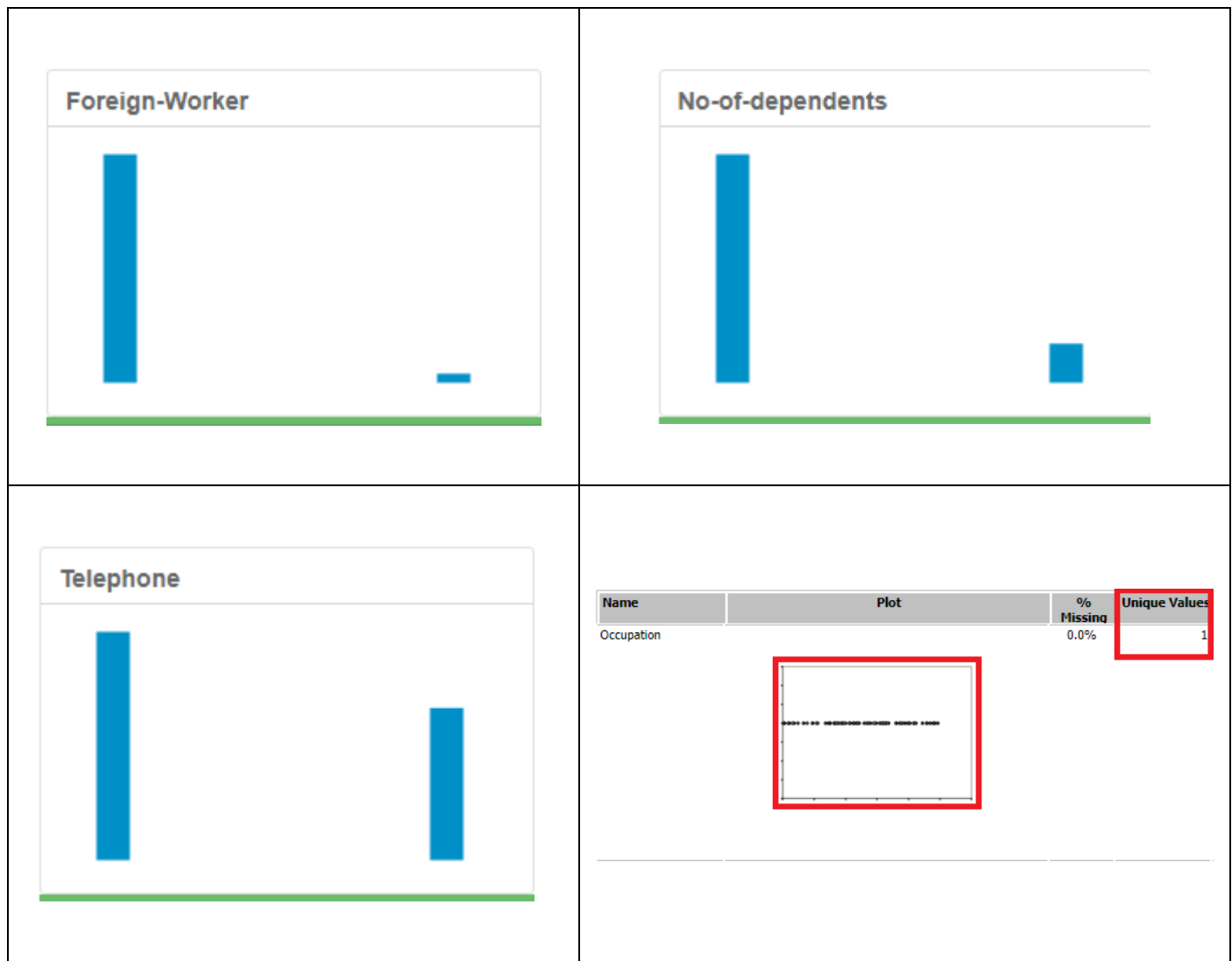


- Are there only a few values in a subset of your data field? Does the data field look very uniform (there is only one value for the entire field?). This is called "low variability" and you should remove fields that have low variability. Refer to the "Tips" section to find examples of data fields with low-variability.

  The detected fields with low variability were:
  - Concurrent-Credits
  - Guarantors
  - Foreign-Worker
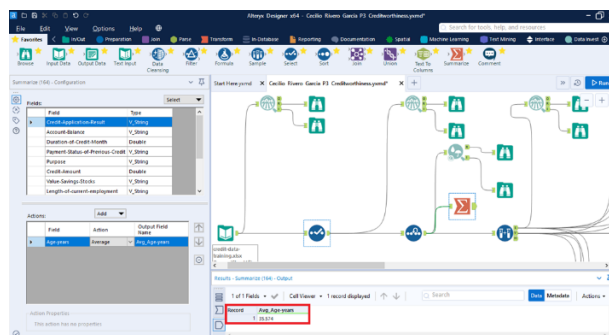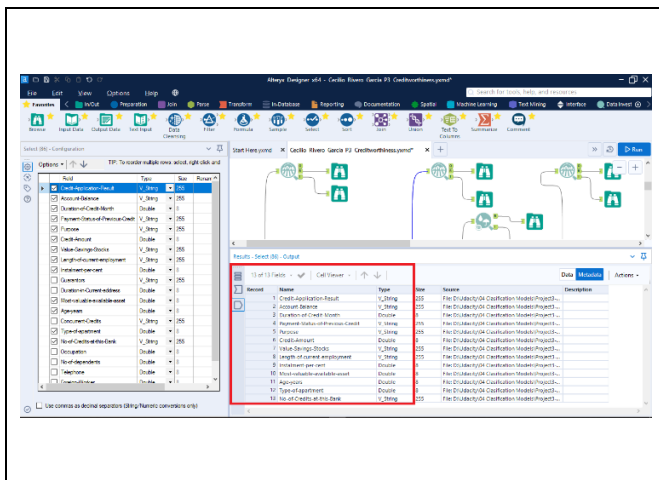  - No-of-dependents
  - Telephone
  - Occupation

**Foreign-Worker**

**No-of-dependents**

**Telephone**

| Name | Plot | % Missing | Unique Values |
|------|------|-----------|---------------|
| Occupation | | 0.0% | 1 |

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

  After analysis of data, the fields chosen to build the predictive models were (13):
  1. Credit-Application-Result
  2. Account-Balance
  3. Duration-of-Credit-Month
  4. Payment-Status-of-Previous-Credit
  5. Purpose
  6. Credit-Amount
  7. Value-Savings-Stocks
  8. Length-of-current-employment
  9. Instalment—per-cent
  10. Most-valuable-available-asset
  11. Age-years ≈ average 36 years
  12. Type-of-apartment
  13. No-of-Credits-at-this-Bank

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

*To achieve consistent results reviewers expect.*

*Answer this question:*

- In your cleanup process, which fields did you remove or impute? Please justify why you removed or imputed these fields. Visualizations are encouraged.
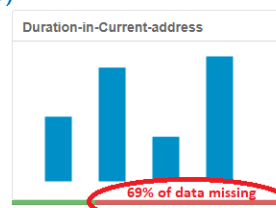  The data fields removed were:

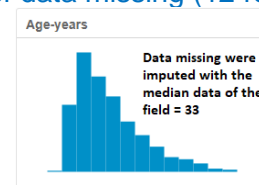| Data Field Removed | Justification |
|---|---|
| ➤ Duration-in-current-address | High amount of missing data (69%)  |
| ➤ Concurrent-Credits<br>➤ Guarantors<br>➤ Foreign-Worker<br>➤ No-of-dependents<br>➤ Telephone<br>➤ Occupation | Low-Variability  |

The data field imputed was:

| Data Field Imputed | Justification |
|---|---|
| ➤ Age-year | 2% of data missing (12 records).  |

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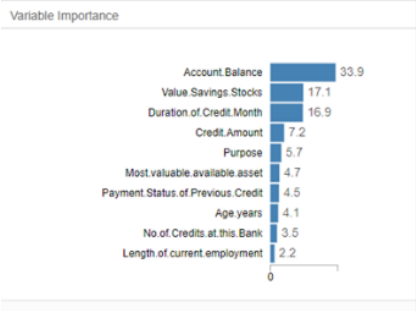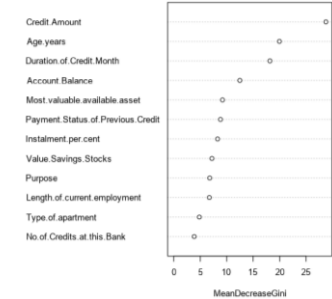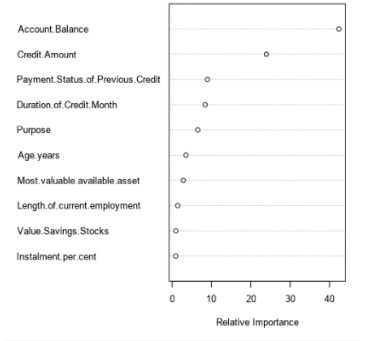
*Answer these questions for **each model** you created:*

- Which predictor variables are significant or the most important? Please show the p-values or variable importance charts for all of your predictor variables.
  **Using Credit-Application-Result as target variable and setup the create sample tool as indicated:**
  - Estimation sample: 70
  - Validation sample: 30
  - Random seed: 1

  The following table show the significant or most important predictor variables by predictive model:

| Predictive Model | Report |
|---|---|
| Logistic Regression | **Most Significant Predictor Variable** — **P-Value** <br><br> Account-Balance/Some Balance — 1.79e-06 <br> Purpose/New Car — 0.00519 <br> Credit-Amount — 0.00989 <br><br>  |
| Logistic Regression-Stepwise | **Most Significant Predictor Variable** — **P-Value** <br><br> Account-Balance/Some Balance — 1.65e-06 <br> Purpose/New Car — 0.00566 <br> Credit-Amount — 0.00296 <br><br>  |

| Predictive Model | Report | |
|---|---|---|
| Decision Tree | **Most Significant Predictor Variable** | **Weighing** |
| | Account-Balance | 33.9 |
| | Value Saving Stocks | 17.1 |
| | Duration Credit Month | 16.9 |



Variable Importance

| Forest Model | Most Significant Predictor Variable | Weighing |
|---|---|---|
| | Credit Amount | >25 |
| | Age years | >20 |
| | Duration of Credit Month | >15 |



Variable Importance Plot

| Boosted Model | Most Significant Predictor Variable | Weighing |
|---|---|---|
| | Account Balance | >40 |
| | Credit Amount | >20 |
| | Payment Status of Previous Credit | >15 |



Variable Importance Plot

- 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?

| Predictive Model | Validation Set |
|---|---|
| Logistic Regression | The overall predictive model accuracy is 78%, and the report shows that it is biased by the "Creditworthy" data with the 90% versus 49% of accuracy by the "Non-Creditworthy" data and a PPV = 81% versus NPV= 69% respectively.<br><br>**Model Comparison Report**<br>**Fit and error measures**<br>Model: LoR, Accuracy 0.7800, F1 0.8520, AUC 0.7314, Accuracy_Creditworthy 0.9045, Accuracy_Non-Creditworthy 0.4889<br>Model: model names in the current comparison.<br>Accuracy: overall accuracy, number of correct predictions of all classes divided by total sample number.<br>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.<br>AUC: area under the ROC curve, only available for two-class classification.<br>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.<br><br>**Confusion matrix of LoR**<br>| | Actual_Creditworthy | Actual_Non-Creditworthy |<br>Predicted_Creditworthy 95, 23<br>Predicted_Non-Creditworthy 10, 22 |
| Logistic Regression-Stepwise | The overall predictive model accuracy is 76%, and the report shows that it is biased by the "Creditworthy" data with the 88% versus 49% of accuracy by the "Non-Creditworthy" data and a PPV = 80% versus NPV= 63% respectively.<br><br>**Model Comparison Report**<br>**Fit and error measures**<br>Model: StepwiseLoR, Accuracy 0.7600, F1 0.8364, AUC 0.7306, Accuracy_Creditworthy 0.8762, Accuracy_Non-Creditworthy 0.4889<br>Model: model names in the current comparison.<br>Accuracy: overall accuracy, number of correct predictions of all classes divided by total sample number.<br>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.<br>AUC: area under the ROC curve, only available for two-class classification.<br>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.<br><br>**Confusion matrix of StepwiseLoR**<br>| | Actual_Creditworthy | Actual_Non-Creditworthy |<br>Predicted_Creditworthy 92, 23<br>Predicted_Non-Creditworthy 13, 22 |
| Decision Tree | The overall predictive model accuracy is 75%, and the report shows that it is biased by the "Creditworthy" data with the 89% versus 42% of accuracy by the "Non-Creditworthy" data and a PPV = 78% versus NPV= 61% respectively.<br><br>**Model Comparison Report**<br>**Fit and error measures**<br>Model: DTM, Accuracy 0.7467, F1 0.8304, AUC 0.7035, Accuracy_Creditworthy 0.8857, Accuracy_Non-Creditworthy 0.4222<br>Model: model names in the current comparison.<br>Accuracy: overall accuracy, number of correct predictions of all classes divided by total sample number.<br>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.<br>AUC: area under the ROC curve, only available for two-class classification.<br>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.<br><br>**Confusion matrix of DTM**<br>| | Actual_Creditworthy | Actual_Non-Creditworthy |<br>Predicted_Creditworthy 93, 26<br>Predicted_Non-Creditworthy 12, 19<br>**Performance Diagnostic Plots** |
| Forest Model | The overall predictive model accuracy is 79%, and the report shows that it is not biased with a PPV of 79% for "Creditworthy" data versus an NPV of 85% for "Non-Creditworthy" data.<br><br>**Model Comparison Report**<br>**Fit and error measures**<br>Model: FM, Accuracy 0.7933, F1 0.8681, AUC 0.7368, Accuracy_Creditworthy 0.9714, Accuracy_Non-Creditworthy 0.3778<br>Model: model names in the current comparison.<br>Accuracy: overall accuracy, number of correct predictions of all classes divided by total sample number.<br>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.<br>AUC: area under the ROC curve, only available for two-class classification.<br>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.<br><br>**Confusion matrix of FM**<br>| | Actual_Creditworthy | Actual_Non-Creditworthy |<br>Predicted_Creditworthy 102, 28<br>Predicted_Non-Creditworthy 3, 17 |
| Boosted Model | The overall predictive model accuracy is 71%, and the report shows that it is not biased with a PPV of 72% for "Creditworthy" data versus an NPV of 67% for "Non-Creditworthy" data.<br><br>**Model Comparison Report**<br>**Fit and error measures**<br>Model: BM, Accuracy 0.7133, F1 0.8273, AUC 0.7368, Accuracy_Creditworthy 0.9810, Accuracy_Non-Creditworthy 0.0889<br>Model: model names in the current comparison.<br>Accuracy: overall accuracy, number of correct predictions of all classes divided by total sample number.<br>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.<br>AUC: area under the ROC curve, only available for two-class classification.<br>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.<br><br>**Confusion matrix of BM**<br>| | Actual_Creditworthy | Actual_Non-Creditworthy |<br>Predicted_Creditworthy 103, 41<br>Predicted_Non-Creditworthy 2, 4<br>**Performance Diagnostic Plots** |

# 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)*

After setup and running all the classification predictive models, the Forest Model was selected, because it shows the highest overall accuracy 79.33%, with an accuracy of 97.14% to predict creditworthy loan requester and also it shows the faster rate to reach true-positives according to the ROC graph, all these give us the security to loan money to the correct applicant with a high probability of recover

Finally applying the Forest Model the result are 408 viable applications (creditworthy) of 500 new loan applications

*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:
    - Overall Accuracy against your Validation set
    - Accuracies within "Creditworthy" and "Non-Creditworthy" segments
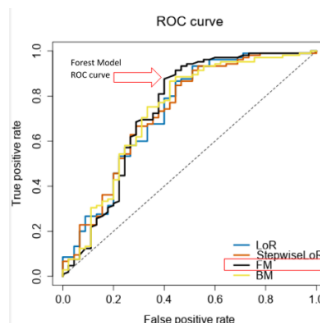    - ROC graph
    - Bias in the Confusion Matrices

The predictive model with the best performance is Forest Model, showing the following parameters:
- Overall accuracy: 79.33%
- Creditworthy accuracy: 97.14%
- Non-Creditworthy accuracy: 37.78%

Layout

## Model Comparison Report

### Fit and error measures

| Model | Accuracy | F1 | AUC | Accuracy_Creditworthy | Accuracy_Non-Creditworthy |
|---|---|---|---|---|---|
| LoR | 0.7800 | 0.8520 | 0.7314 | 0.9048 | 0.4889 |
| StepwiseLoR | 0.7600 | 0.8364 | 0.7306 | 0.8762 | 0.4889 |
| FM | 0.7933 | 0.8681 | 0.7368 | 0.9714 | 0.3778 |
| BM | 0.7133 | 0.8273 | 0.7308 | 0.9810 | 0.0889 |

**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.

The ROG graphic also shows the best performance related with true-positive rate of prediction.



ROC curve

| Confusion matrix of BM | | |
| --- | --- | --- |
| | Actual_Creditworthy | Actual_Non-Creditworthy |
| Predicted_Creditworthy | 103 | 41 |
| Predicted_Non-Creditworthy | 2 | 4 |

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

| Confusion matrix of LoR | | |
| --- | --- | --- |
| | Actual_Creditworthy | Actual_Non-Creditworthy |
| Predicted_Creditworthy | 95 | 23 |
| Predicted_Non-Creditworthy | 10 | 22 |

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

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

- How many individuals are creditworthy?
  408 are the new customers (individuals) that according to the predictive Forest- Model will be creditworthy to pay the loan requested.

Alteryx's Workflows:

1. Workflow 1(Model Comparative & Final Predictive Modeler):



2. Workflow 2(Validation of Final Predictive Modeler):