# Introduction to Data Analytics

Assignment 3                                                      Kaung Khant Kyaw
                                                                        14447837

# Data Mining Problem

This assignment mainly focuses on data training and model optimization. This assignment required me to predict previous data by using and modifying classifiers as a data scientist in a company. I used KNIME software as a tool to see the dataset and use models in the tool to finish this assignment.

## Input

The dataset given is about National Basketball Association players and their stats to train a model that can predict if they will be playing in the next 5 years future.
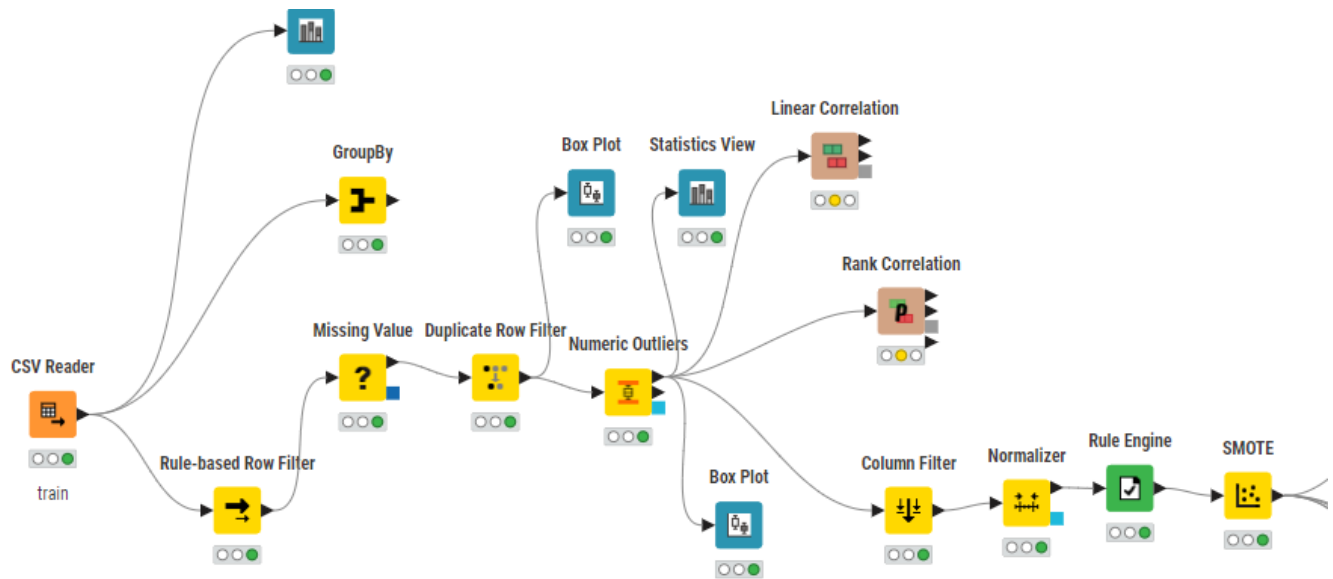
## Output

The problem is that the data given includes some negative values as well as outliers, and since negative values are unstable to predict a model, we would have to pre-process this data first before training a model to predict. We also had to remove the outliers too.

## Goal

The goal of this assessment is to demonstrate the ability to preprocess a dataset, train a predictive model, and optimize the classifier for predicting future outcomes. Identify and handle negative values in the dataset, ensuring the data is suitable for training a predictive model. Detect and remove outliers that could skew the model's performance. Utilize KNIME software to train a model on the preprocessed dataset. Experiment with various classifiers to determine the best fit for the data. Modify and fine-tune the selected classifier(s) to improve prediction accuracy. Evaluate the performance of the model using appropriate metrics. Use the trained and optimized model to predict whether NBA players will be playing in the next 5 years based on their stats.
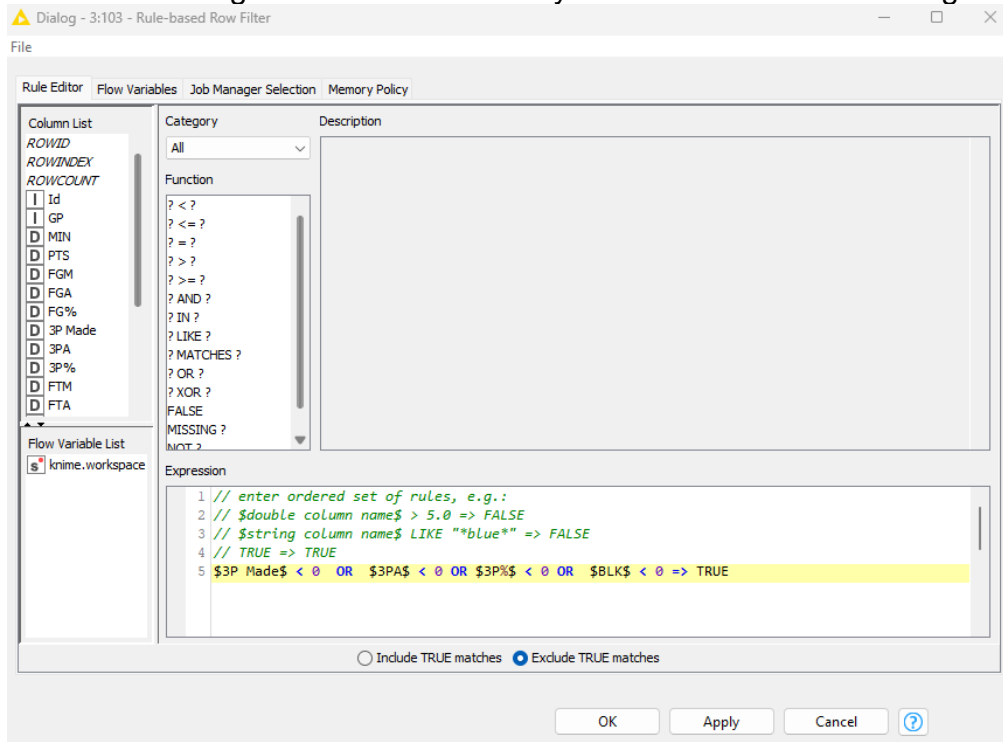
# Data Pre-Processing



## CSV reader

Input the dataset given as a csv file type which is train.csv.
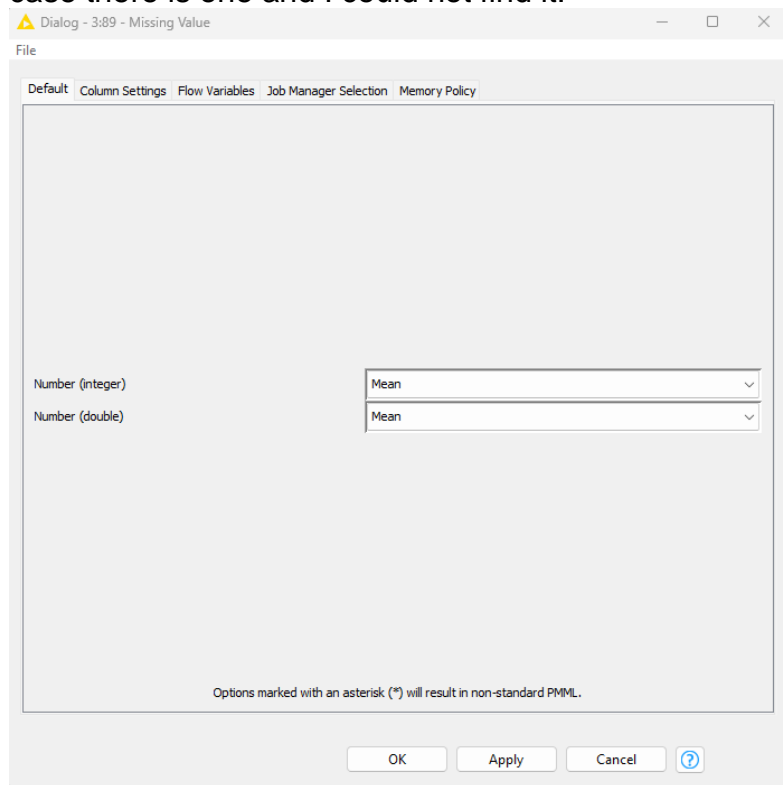
## Rule-based row Filter
Filter out the negative values since they are unstable when building a model.

## Missing Value

Even though, there were not any missing values, I set the missing values to the mean values just in case there is one and I could not find it.
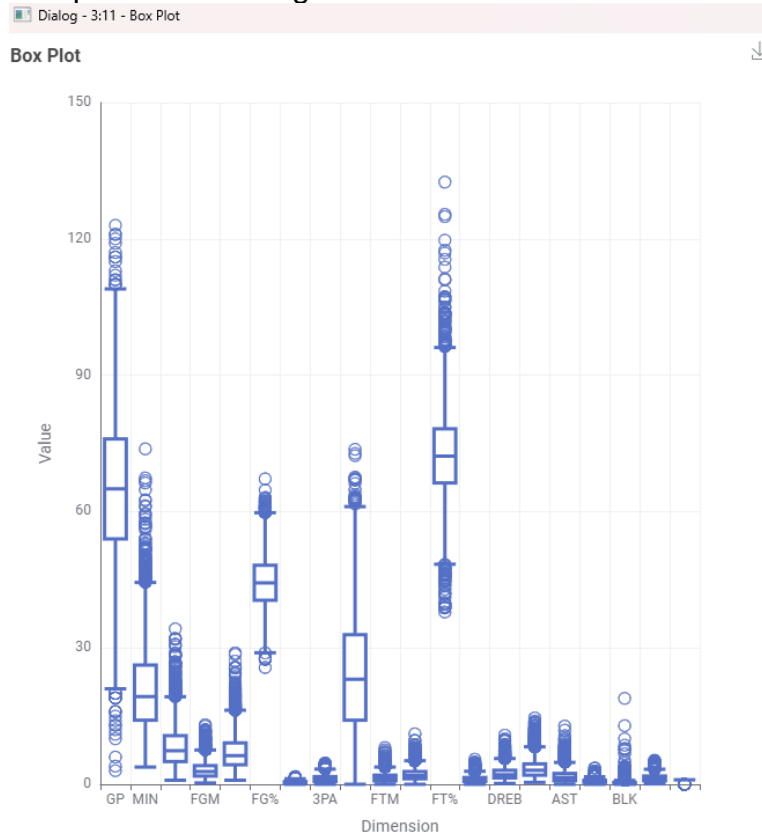


## Duplicate row filter

This node is for checking if there were any duplicate rows in the dataset, luckily there were not any, so I did not have to deal with it.

# Numeric Outliers

I also used box plot to check the outliers before and after using numeric outliers which is used to replace the outlier values.

Box plot before using this node



## Numeric Outliers
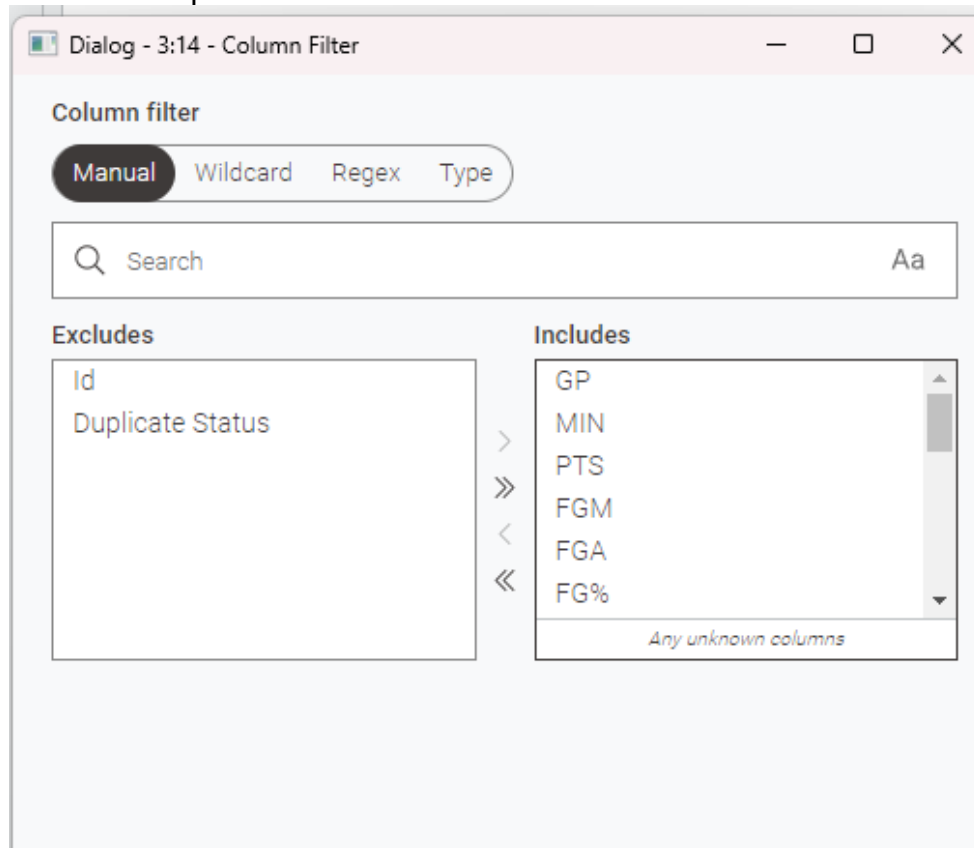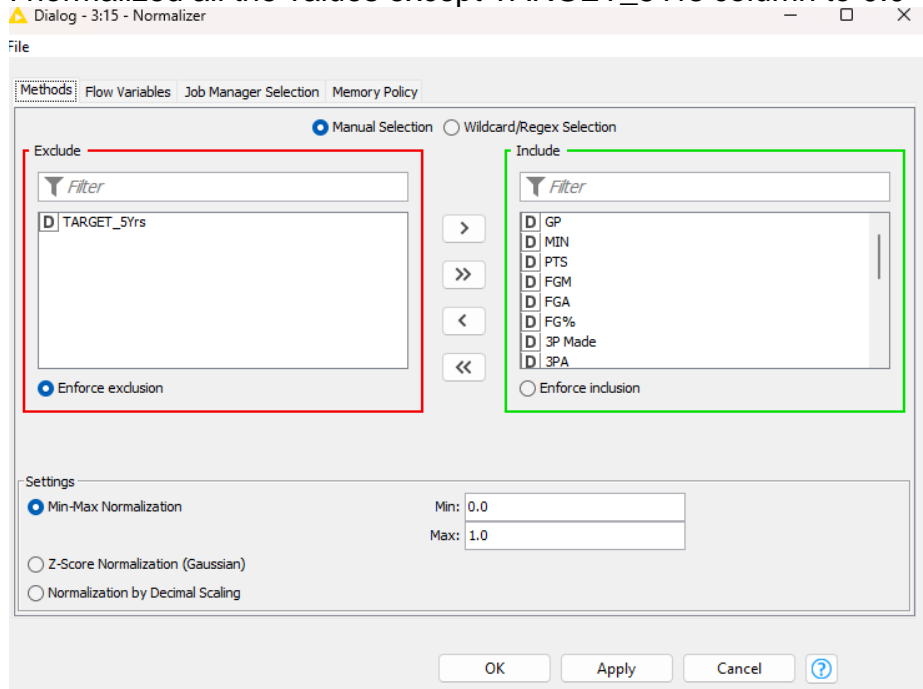
After using Numeric Outliers



## Column Filter

In the column filter, I filtered out the player's ID. And the new duplicate status that is used to check if the row is duplicated or not.
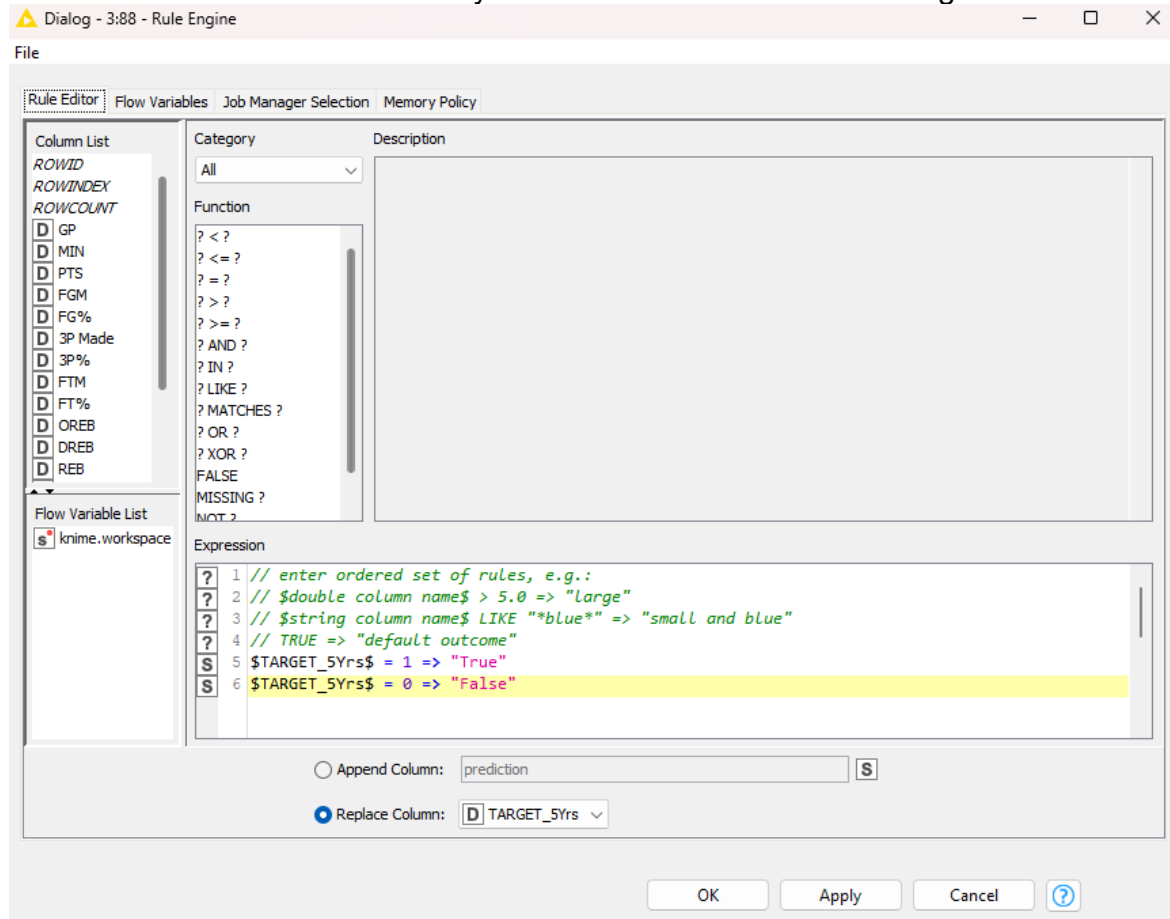
# Normalizer

I normalized all the values except TARGET_5Yrs column to 0.0-1.0 ratio.



# Rule engine

I used rule-engine to change my TARGET_5Yrs values (1,0) to (True,False) and nominal data type for the client to understand clearly and nominal for future modelling.

## Smote

SMOTE (Synthetic Minority Over-sampling Technique) is used to oversample the input data to enrich the training data.



## Partitioning

All the partitioning done in this model are 70 percent relative.

## Linear and Rank Correlation



Correlation Matrix - 3:90 - Rank Correlation

File   View

Rows: d, GP, MIN, PTS, FGM, FGA, FG%, 3P Made, 3PA, 3P%, FTM, FTA, FT%, OREB, DREB, REB, AST, STL, BLK, TOV, TARGET_5Yrs

This matrix provides a quick overview of how pairs of variables are related, helping identify important relationships that might warrant further investigation or consideration in predictive models

# Classifiers & Explanations

## All predictors

I append a new column with a normalized class distribution to be able to see newly predicted columns in the ROC curve. This applies to all predictors used in this assessment.
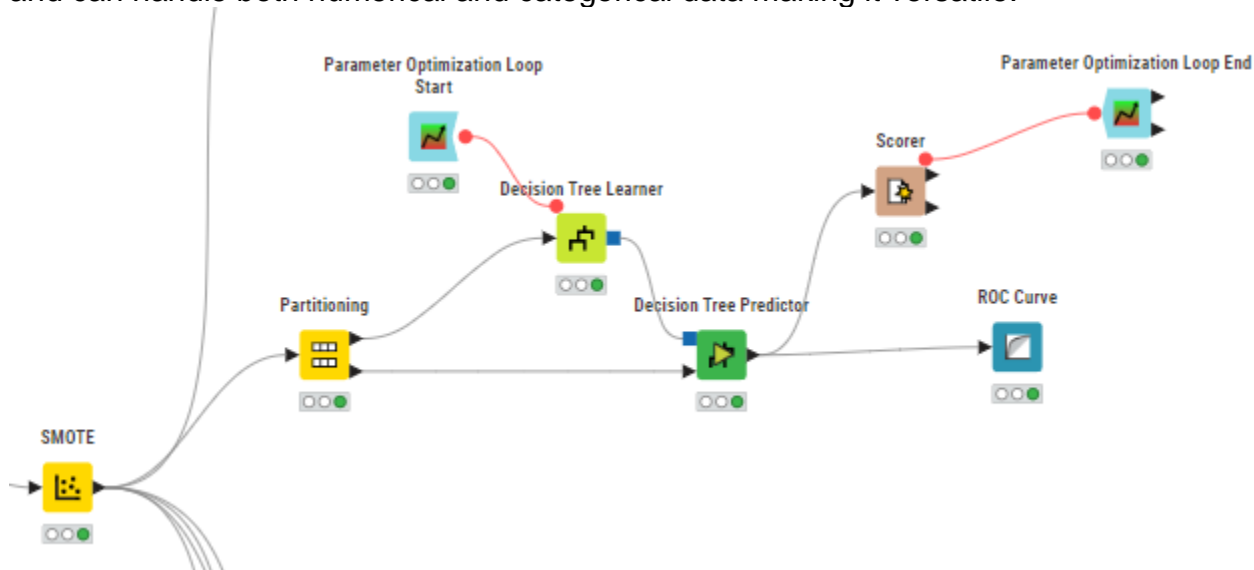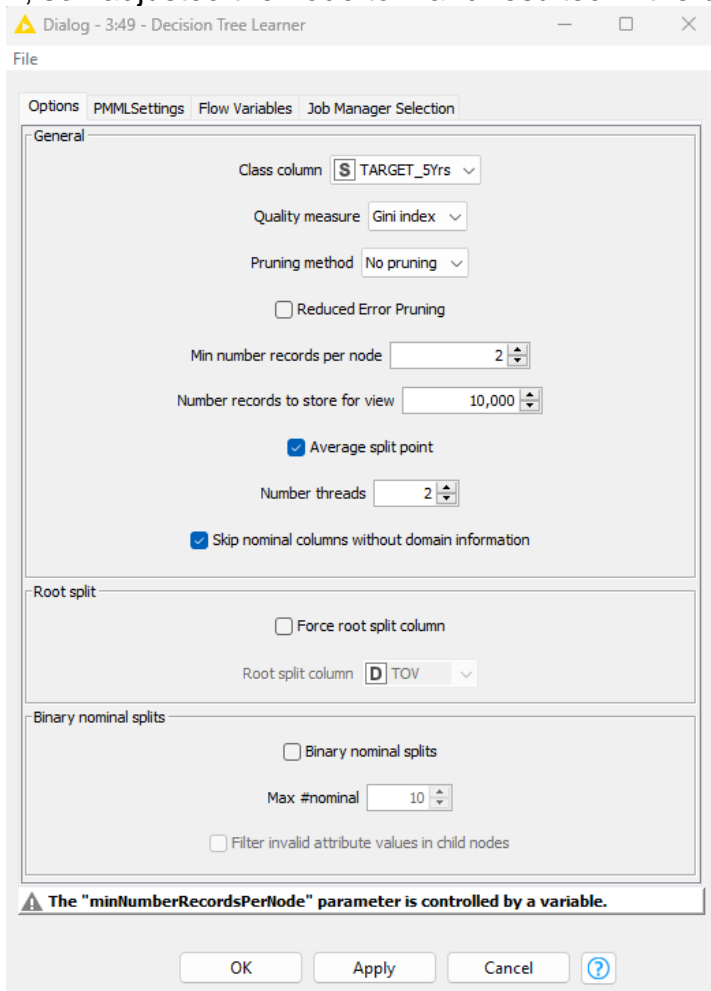


## Decision Tree

This classifier worked and the data is predicted into structure trees. This is effective for noisy data and can handle both numerical and categorical data making it versatile.



I used parameter optimization loop to adjust the parameter for the learner to have the highest accuracy for the score.

The loop shows that the accuracy would be the highest if the minimum number of records per node is 2, so I adjusted the node to 2 and resulted in the best accuracy possible for this classifier.

# Here is the accuracy of the classifier



| Prediction ... | True | False |
|---|---|---|
| True | 3577 | 318 |
| False | 214 | 350 |

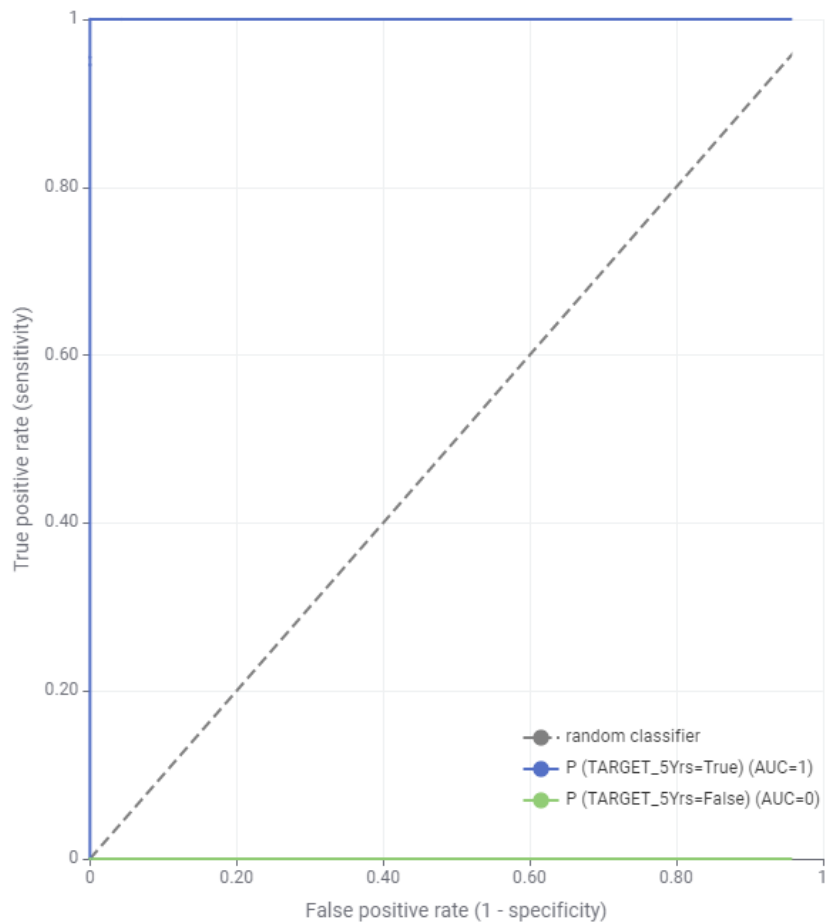Correct classified: 3,927      Wrong classified: 532

Accuracy: 88.069%      Error: 11.931%

Cohen's kappa (κ): 0.5%

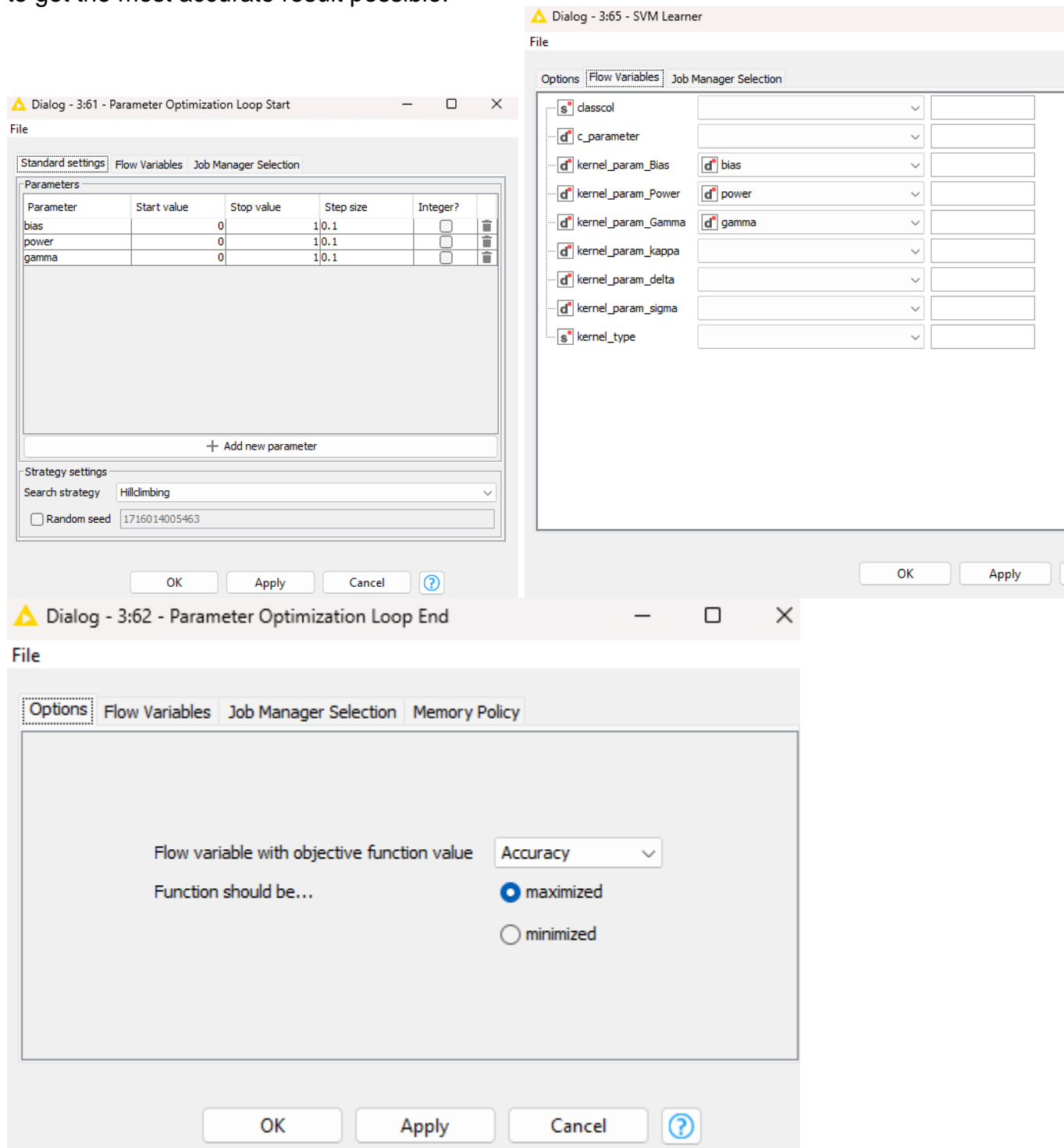# ROC Curve

## SVM Learner

Just as Decision Tree Learner, I used parameter optimization loop to adjust the polynomial variables to get the most accurate result possible.



Using this finds the variables input for the best accuracy and, the settings of SVM learner are

The score of this classifier is



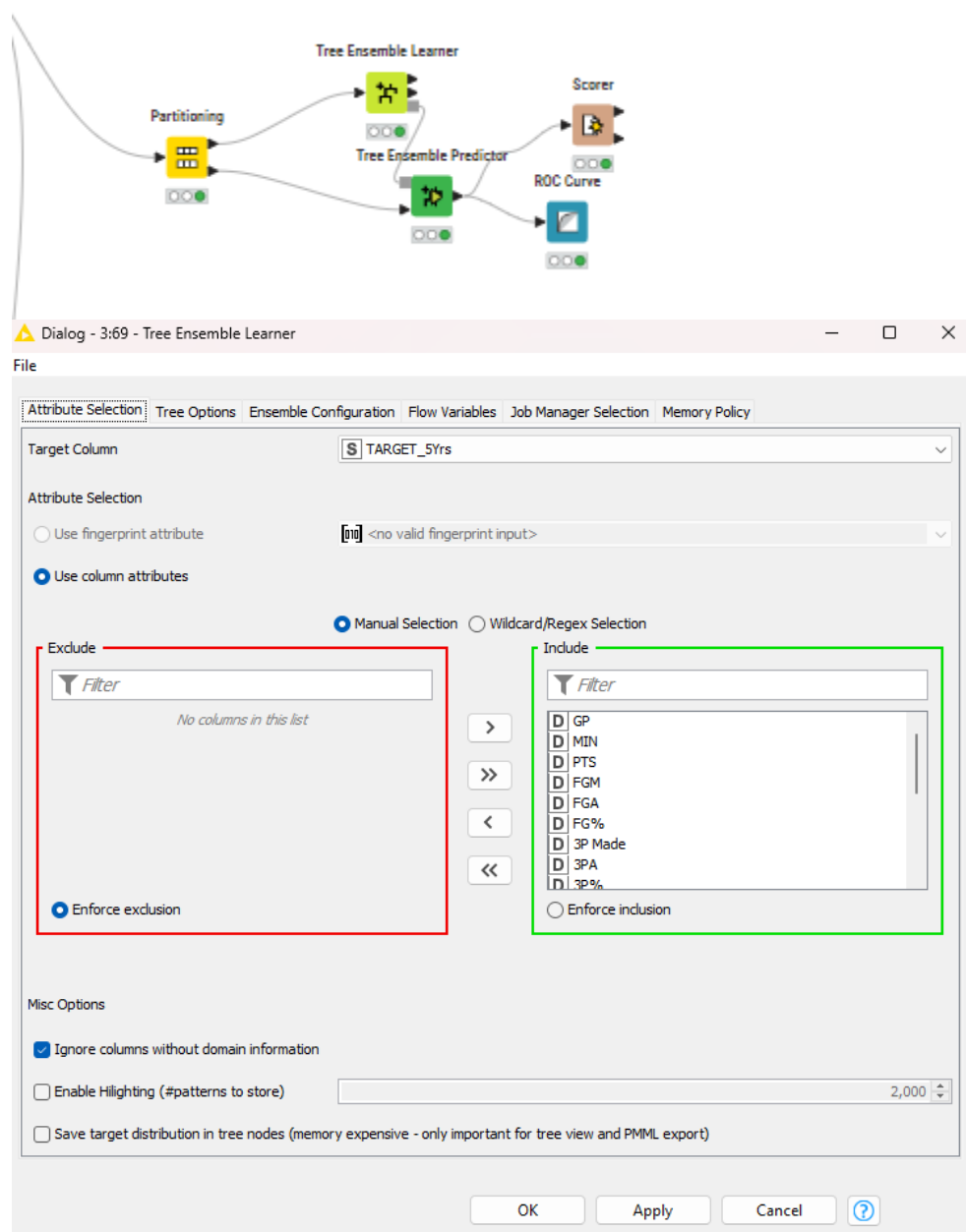| TARGET_5... | True | False |
|---|---|---|
| True | 3758 | 0 |
| False | 701 | 0 |

Correct classified: 3,758    Wrong classified: 701
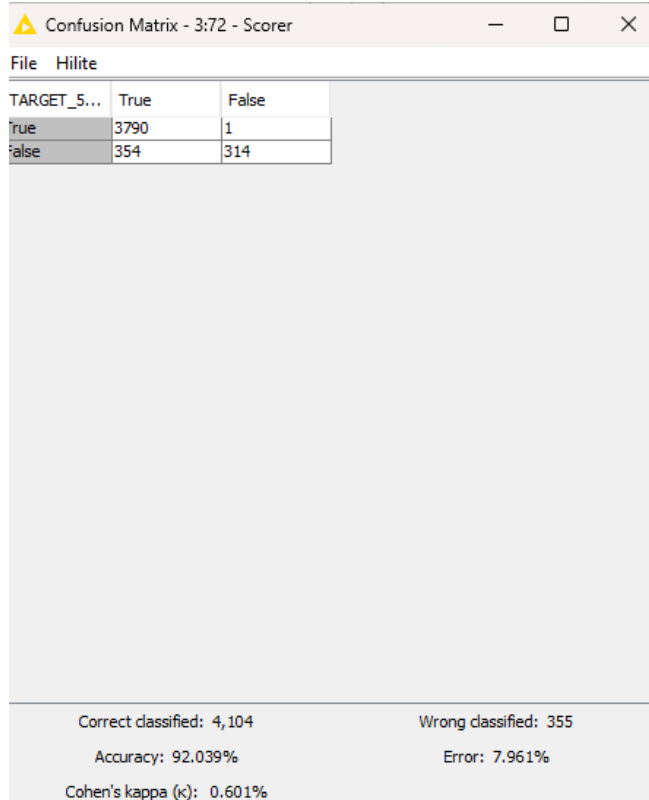
Accuracy: 84.279%    Error: 15.721%
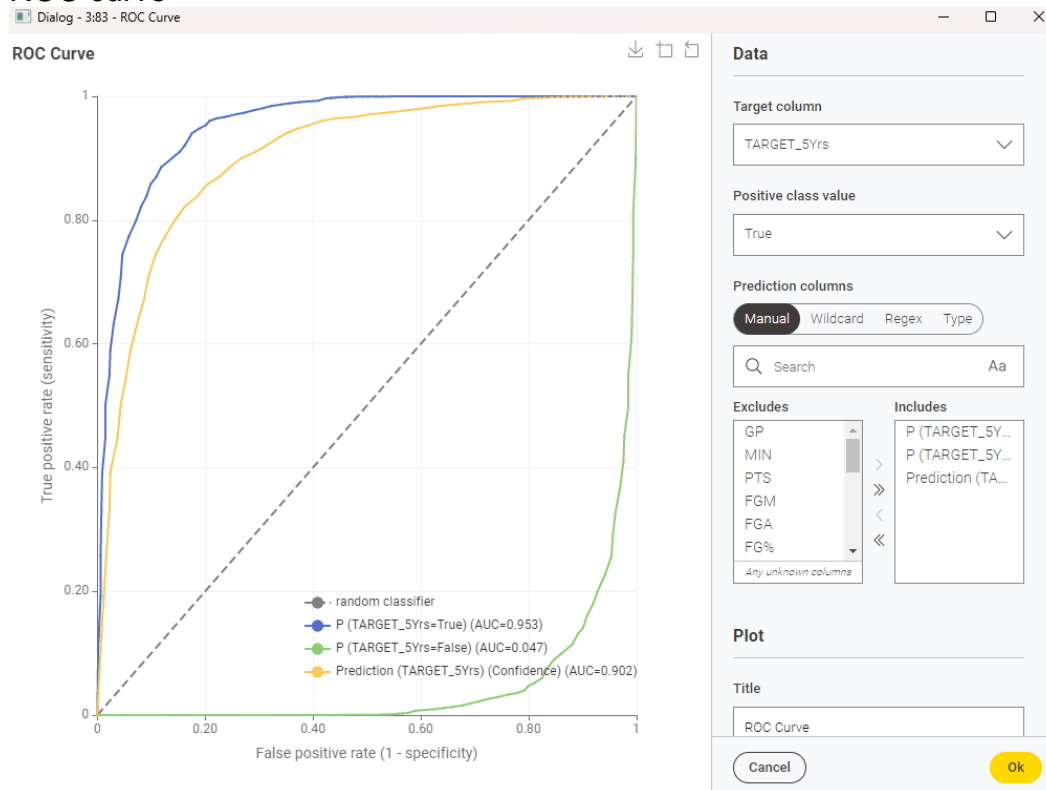
Cohen's kappa (κ): 0%

# Tree ensemble Learner

As there are no changeable parameters in Tree ensemble learner, I used the default settings for the learner.
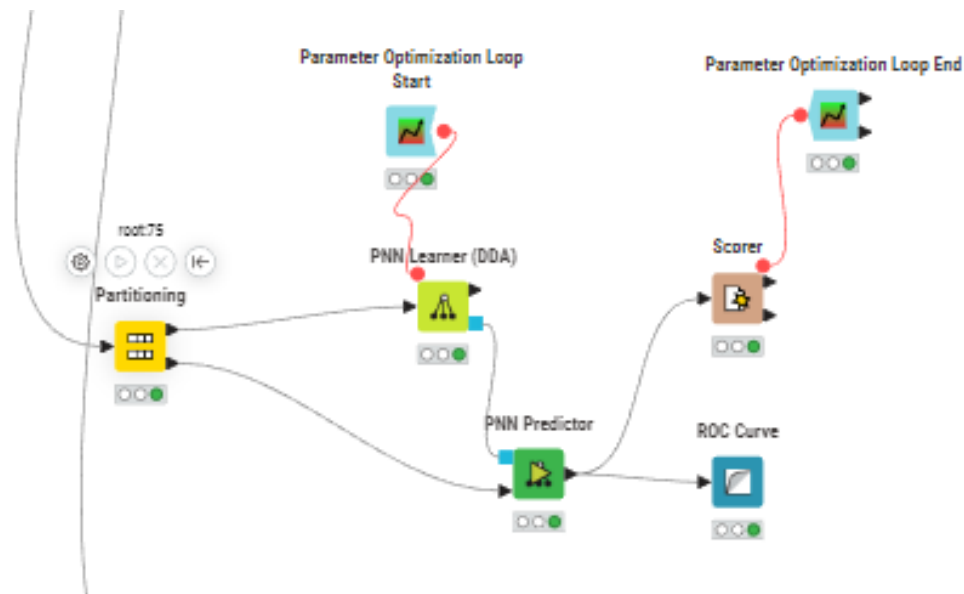
# The score of this classifier is



Confusion Matrix - 3:72 - Scorer

File  Hilite

| TARGET_5... | True | False |
|---|---|---|
| True | 3790 | 1 |
| False | 354 | 314 |

Correct classified: 4,104        Wrong classified: 355

Accuracy: 92.039%        Error: 7.961%
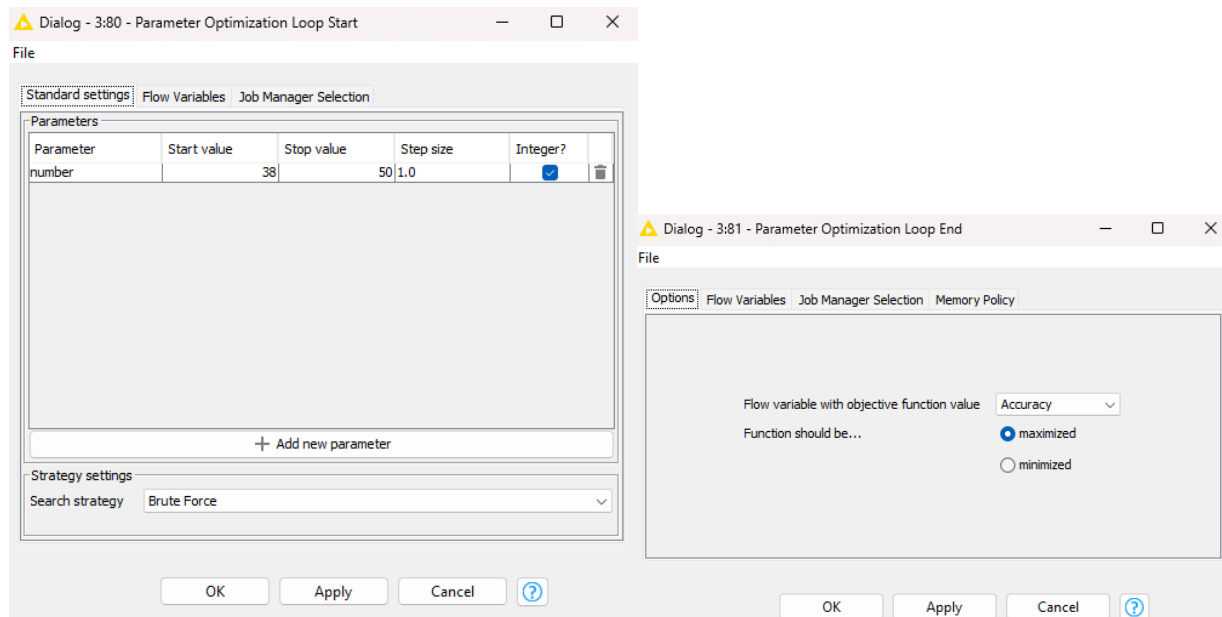
Cohen's kappa (κ): 0.601%

# ROC curve

# PNN learner



In the PNN learner, I used parameter optimization loop to find out the best possible accuracy, I made a variable for the maximum no of epochs. And the other settings are default

The score of the classifier is as follows

## Confusion Matrix - 3:78 - Scorer

File   Hilite

| TARGET_5... | True | False |
|-------------|------|-------|
| True        | 3790 | 1     |
| False       | 638  | 30    |

Correct classified: 3,820          Wrong classified: 639

Accuracy: 85.669%          Error: 14.331%

Cohen's kappa (κ): 0.074%

## ROC curve

### ROC Curve

Dialog - 3:82 - ROC Curve

**Data**

Target column

TARGET_5Yrs

Positive class value

True

Prediction columns

Manual   Wildcard   Regex   Type

Search   Aa

Excludes

DREB
REB
AST
STL
BLK
TOV
Any unknown columns

Includes

P (TARGET_5Y...
P (TARGET_5Y...

**Plot**

Title

ROC Curve

Cancel          Ok

- random classifier
- P (TARGET_5Yrs=False) (AUC=0.167)
- P (TARGET_5Yrs=True) (AUC=0.833)

True positive rate (sensitivity)

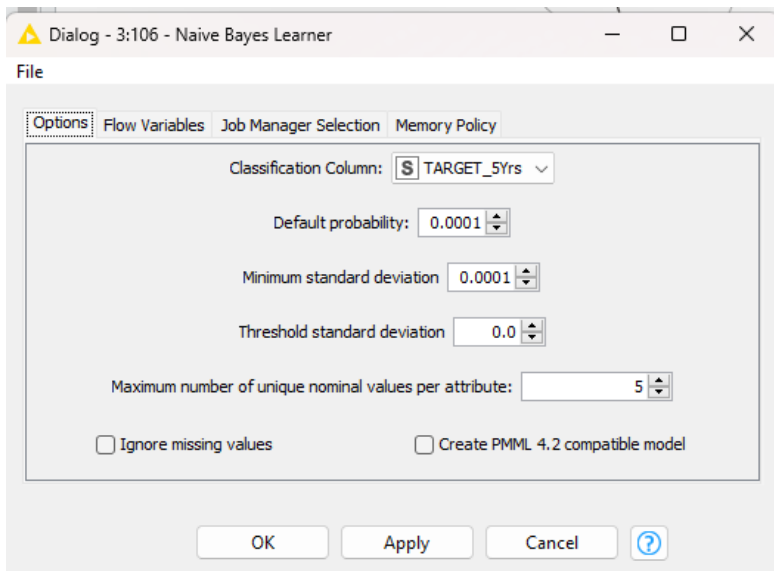False positive rate (1 - specificity)

# Naive Bayes Learner



Naive Bayes Learner with the default settings



The score of Naive Bayes Learner

Confusion Matrix - 3:109 - Scorer

File   Hilite

| TARGET_5... | True | False |
|---|---|---|
| True | 2366 | 1425 |
| False | 191 | 477 |

Correct classified: 2,843          Wrong classified: 1,616

Accuracy: 63.759%                  Error: 36.241%

Cohen's kappa (κ):  0.192%
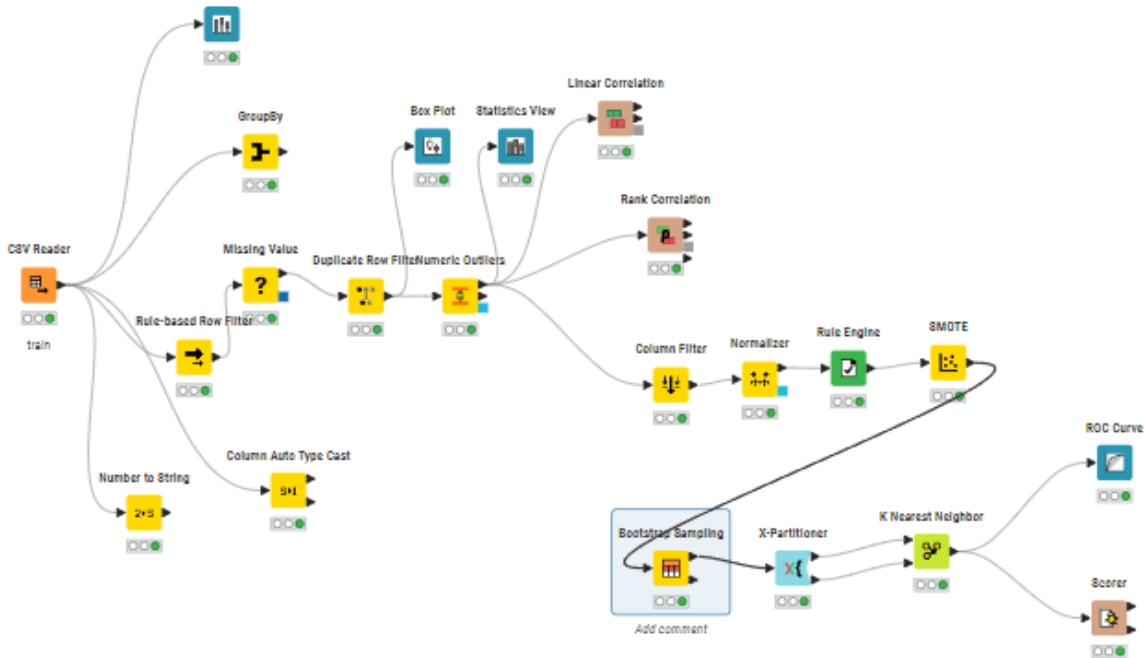
ROC curve

# KNN

In the KNN classifier, I used Bootstrap Sampling node before partitioning, and X- partitioner for the partitioning node.

IN KNN, I used k value 4 so that it gives the higher accuracy than the default settings.



Here is the score of the classifier
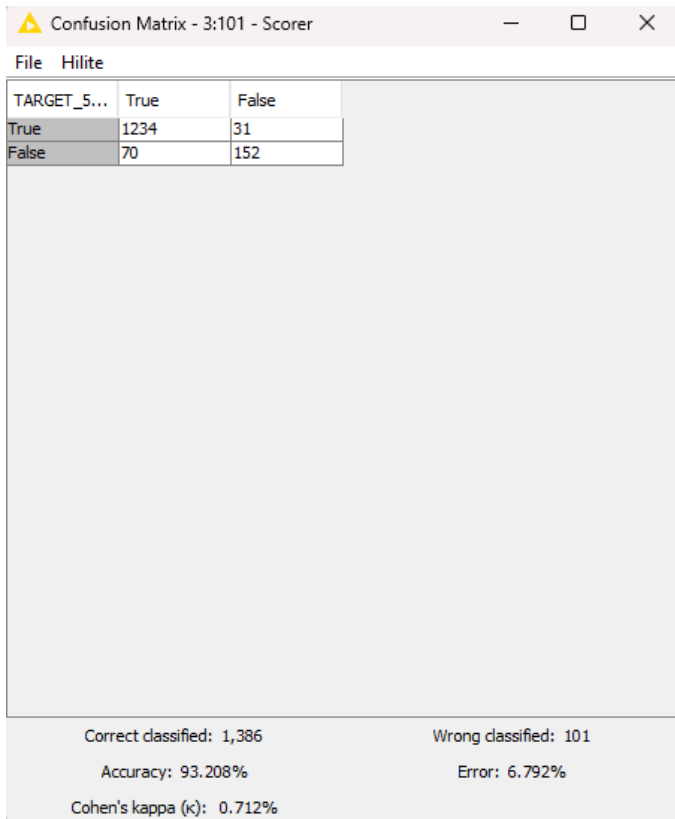
**Confusion Matrix - 3:101 - Scorer**

File   Hilite

| TARGET_5... | True | False |
|---|---|---|
| True | 1234 | 31 |
| False | 70 | 152 |

Correct classified: 1,386          Wrong classified: 101

Accuracy: 93.208%                    Error: 6.792%

Cohen's kappa (κ): 0.712%

## ROC curve

# Random Forest Learner



I used random forest learner default settings



The score of this classifier is

Confusion Matrix - 3:115 - Scorer — □ ✕

File   Hilite

| Prediction ... | True | False |
|---|---|---|
| True | 3787 | 350 |
| False | 4 | 318 |

Correct classified: 4,105                Wrong classified: 354

Accuracy: 92.061%                          Error: 7.939%

Cohen's kappa (κ): 0.604%

ROC Curve

# Best Classifier

| | Accuracy | Recall | Precision | F1 score |
|---|---|---|---|---|

| | | | |
|---|---|---|---|
| SVM | 84.279 | 1 | 1 | |
| Decision Tree | 88.069 | 0.9183 | 0.9435 | 0.9307 |
| PNN | 85.669 | 0.8559 | 0.9997 | 0.9222 |
| Naive Bayes | 63.759 | 0.9253 | 0.6241 | 0.7454 |
| KNN | 93.208 | 0.9463 | 0.9754 | 0.9606 |
| Tree Ensemble | 92.039 | 0.9145 | 0.9997 | 0.9952 |
| Random Forest | 92.061 | 0.9989 | 0.9153 | 0.9552 |

With the metrics resulted, its accuracy, recall and precision is high on Tree Ensemble, also F1 score of Tree Ensemble is high.

Random Forest builds multiple decision trees (bagging) and combines their predictions. Each tree is trained on a random subset of features and instances.

Random Forest is robust, handles noisy data, and avoids overfitting. Its high precision suggests it's good at minimizing false positives.