# kdb+/q AutoML Procedure Report

This report outlines the results for a classification problem achieved through running kdb+/q AutoML.

This run started on 2024.09.06 at 19:26:35.159.

## Description of Input Data

The following is a breakdown of information for each of the relevant columns in the dataset:

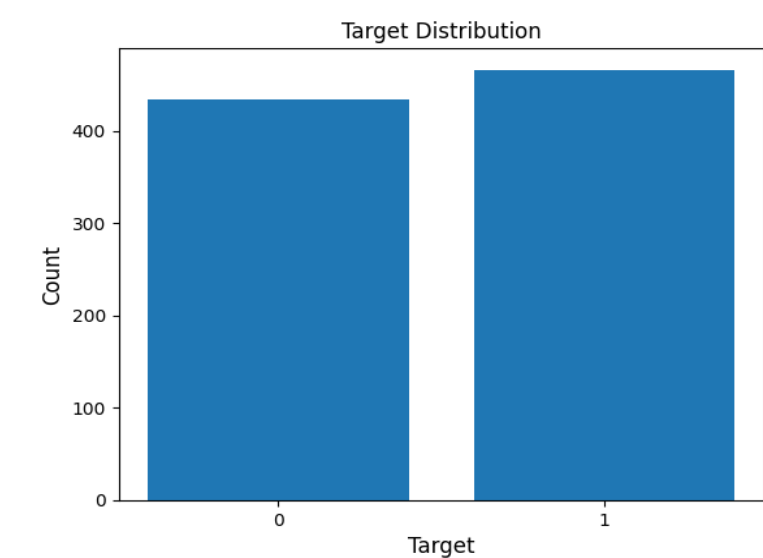| col | count | unique | mean | std | min | max | type |
|-----|-------|--------|------|-----|-----|-----|------|
| comment | 900 | 900 | | | | | text |



Figure 1: Distribution of input target data

## Breakdown of Pre-Processing

Nlp feature extraction and selection was performed with a total of 36 features produced.

Feature extraction took 00:01:12.365 time in total.
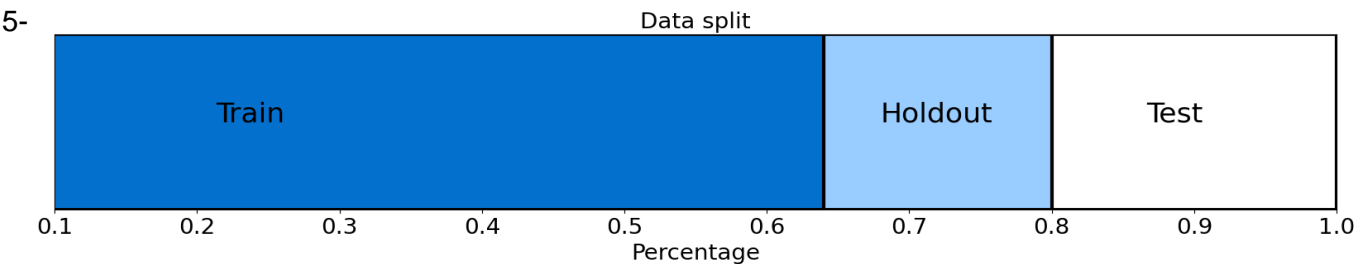
## Initial Scores

Figure 2: The data split used within this run of AutoML, with data split into training, holdout and testing sets

The total time taken to carry out cross validation for each model on the training set was 00:00:07.192 where models were scored and optimized using .ml.accuracy.

Model scores:

RandomForestClassifier = 0.7761469
AdaBoostClassifier = 0.7568966
KNeighborsClassifier = 0.7518291
MLPClassifier = 0.75003
GradientBoostingClassifier = 0.749985
LinearSVC = 0.74997
SVC = 0.7308696
LogisticRegression = 0.7308396
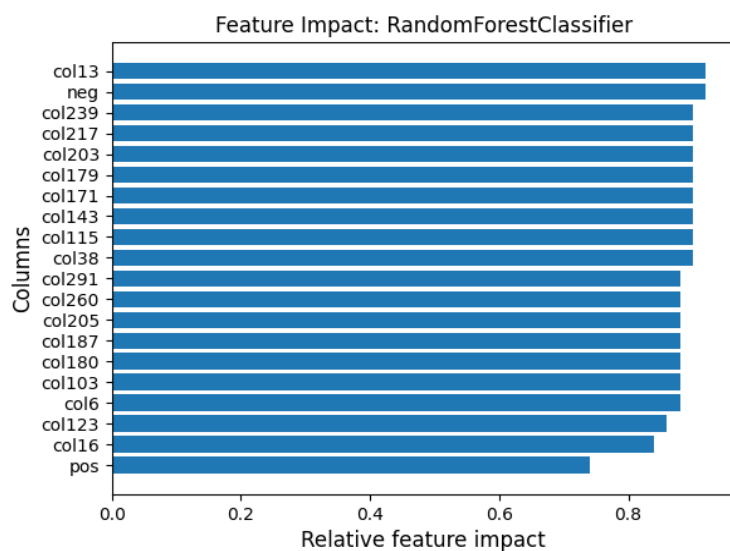BinaryKeras = 0.6926387
GaussianNB = 0.6753823



Figure 3: Feature impact of each significant feature as determined by the training set

## Model selection summary

Best scoring model = RandomForestClassifier

The score on the holdout set for this model was = 0.7430556.

The total time taken to complete the running of this model on the holdout set was: 00:00:00.295.

## Best Model

A 5-fold grid search was performed on the training set to find the best model using, .automl.gs.kfShuff.

The following are the hyperparameters which have been deemed optimal for the model:

criterion = gini
min_samples_split = 2
min_samples_leaf = 1

The score for the best model fit on the entire training set and scored on the testing set was = 0.75
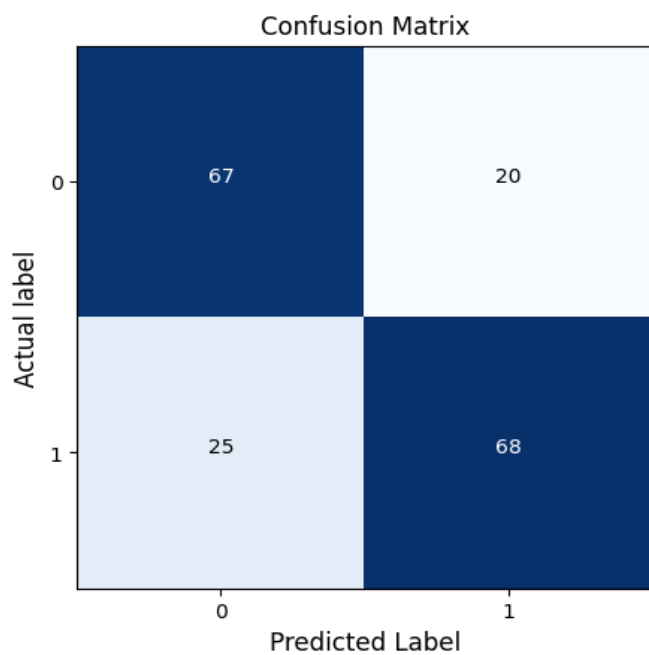


Figure 4: This is the confusion matrix produced for predictions made on the testing set