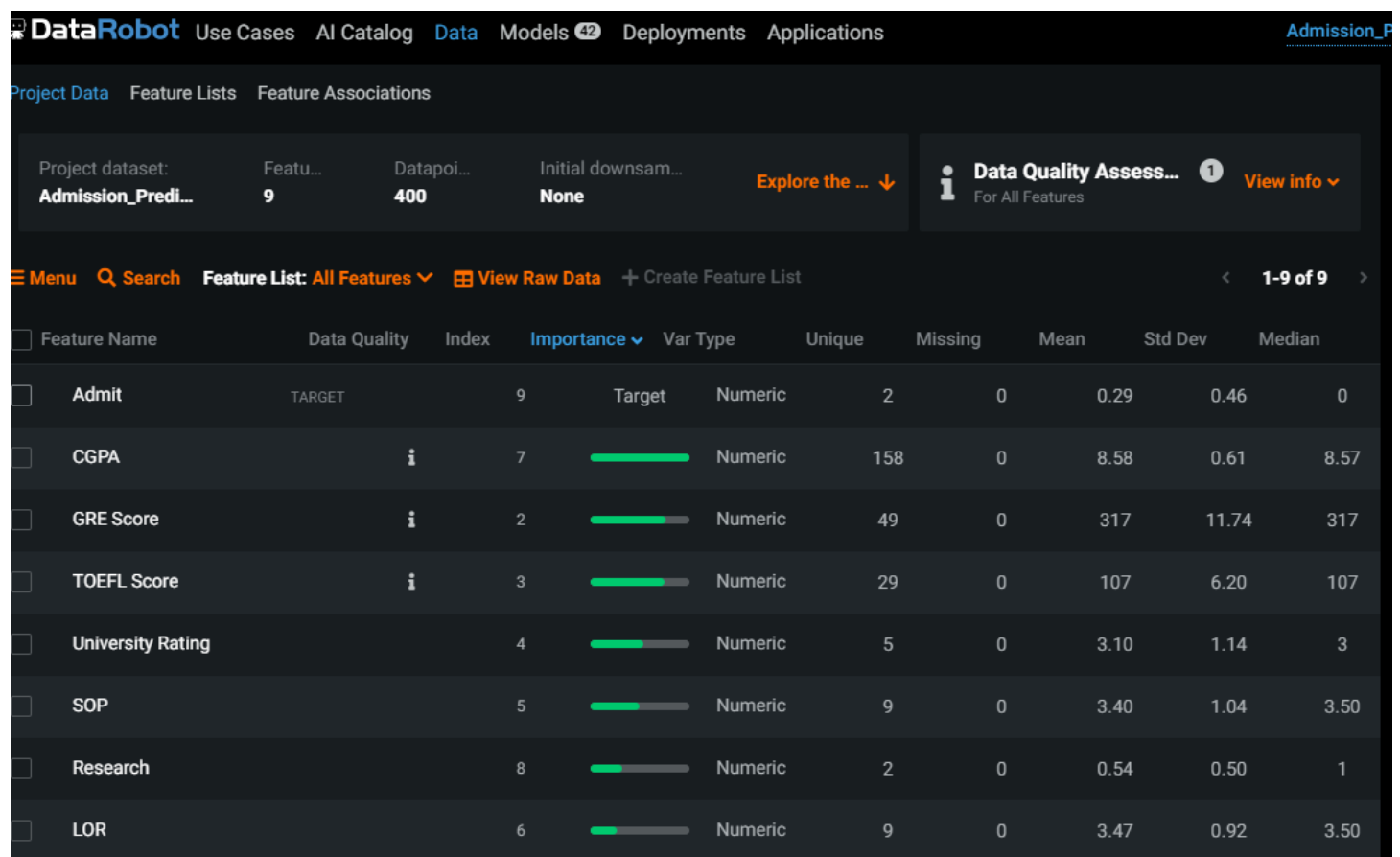**Predictive Modeling with DataRobot**

**by Qi Sun**

The purpose of this study is to predict Graduate Admissions. The dataset contains several variables which are considered important during the application for Masters Programs.

The independent variables are:

GRE Scores ( out of 340 ), TOEFL Scores ( out of 120 ), University Rating ( out of 5 ), Statement of Purpose and Letter of Recommendation Strength ( out of 5 ), Undergraduate GPA ( out of 10 ), Research Experience ( either 0 or 1 ).

The dependent variable is Admit ( 0 or 1 ). There are a total of 400 records in the dataset.

Here is the description of the dataset:



The classification prediction was performed by using DataRobot. I selected 'Autopilot' on the Modeling Mode. After the process completed, I got the results with LogLoss metric showing below:

From the Model tab, we can see that DataRobot has automatically created 41 models. The recommended model is Gradient Boosted Greedy Trees Classifier with 100% data. I decided to pick this Classifier since it has the highest AUC, F1, and accuracy scores and to improve it with hypertuning. Next, I chose the Gradient Boosted Greedy Trees Classifier with 80% training set. The following screenshot shows a description of the model and the evaluation results by using metric of AUC.
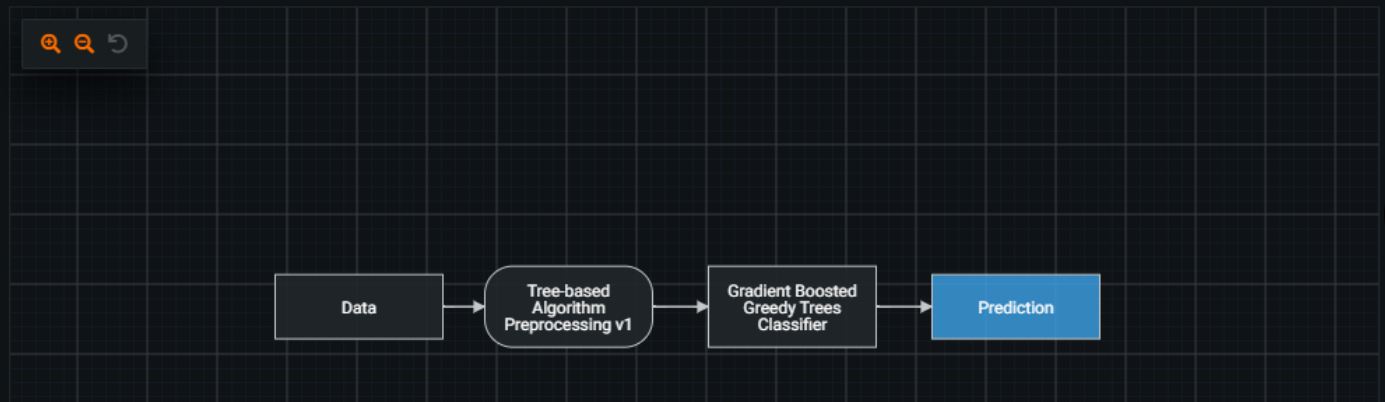
| ☐ Model Name & Description | Feature List & Sample Size ▼ | Validation | Cross Validation | Holdout |
|---|---|---|---|---|

🐍 **Gradient Boosted Greedy Trees Classifier**
Tree-based Algorithm Preprocessing v1

M94  BP75  ❄ 80.0%

Informative Features  ✨  80.0% ➕    0.9450*    0.9706*    0.9794

Evaluate  Understand  **Describe**  Predict  Comments

**Blueprint**  Model Info  Coefficients  Rating Table  Log

| ☐ Model Name & Description | Feature List & Sample Size ▼ | Validation | Cross Validation | Holdout |
|---|---|---|---|---|

🐍 **Gradient Boosted Greedy Trees Classifier**
Tree-based Algorithm Preprocessing v1

M94  BP75  ❄ 80.0%

Informative Features  ✨  80.0% ➕    0.9450*    0.9706*    0.9794

Evaluate  Understand  **Describe**  Predict  Comments

Blueprint  **Model Info**  Coefficients  Rating Table  Log

## Model Overview

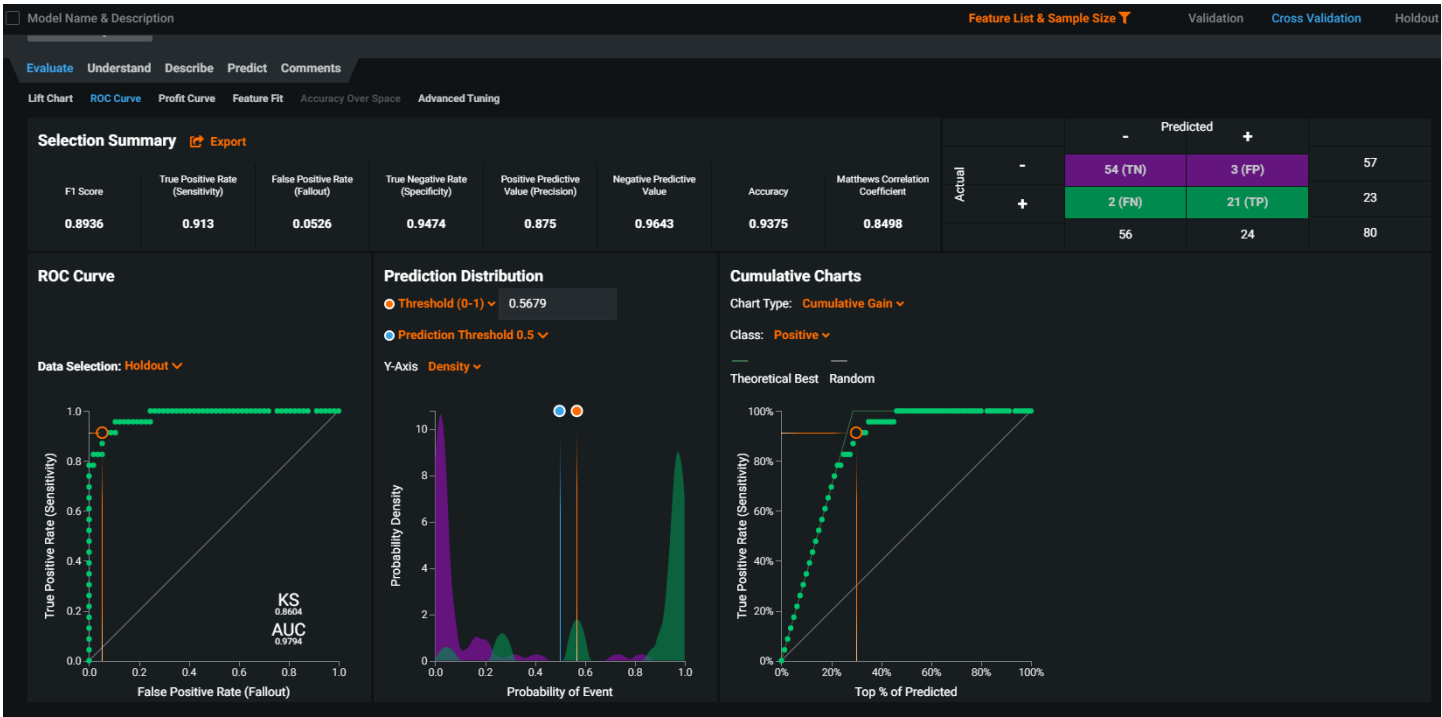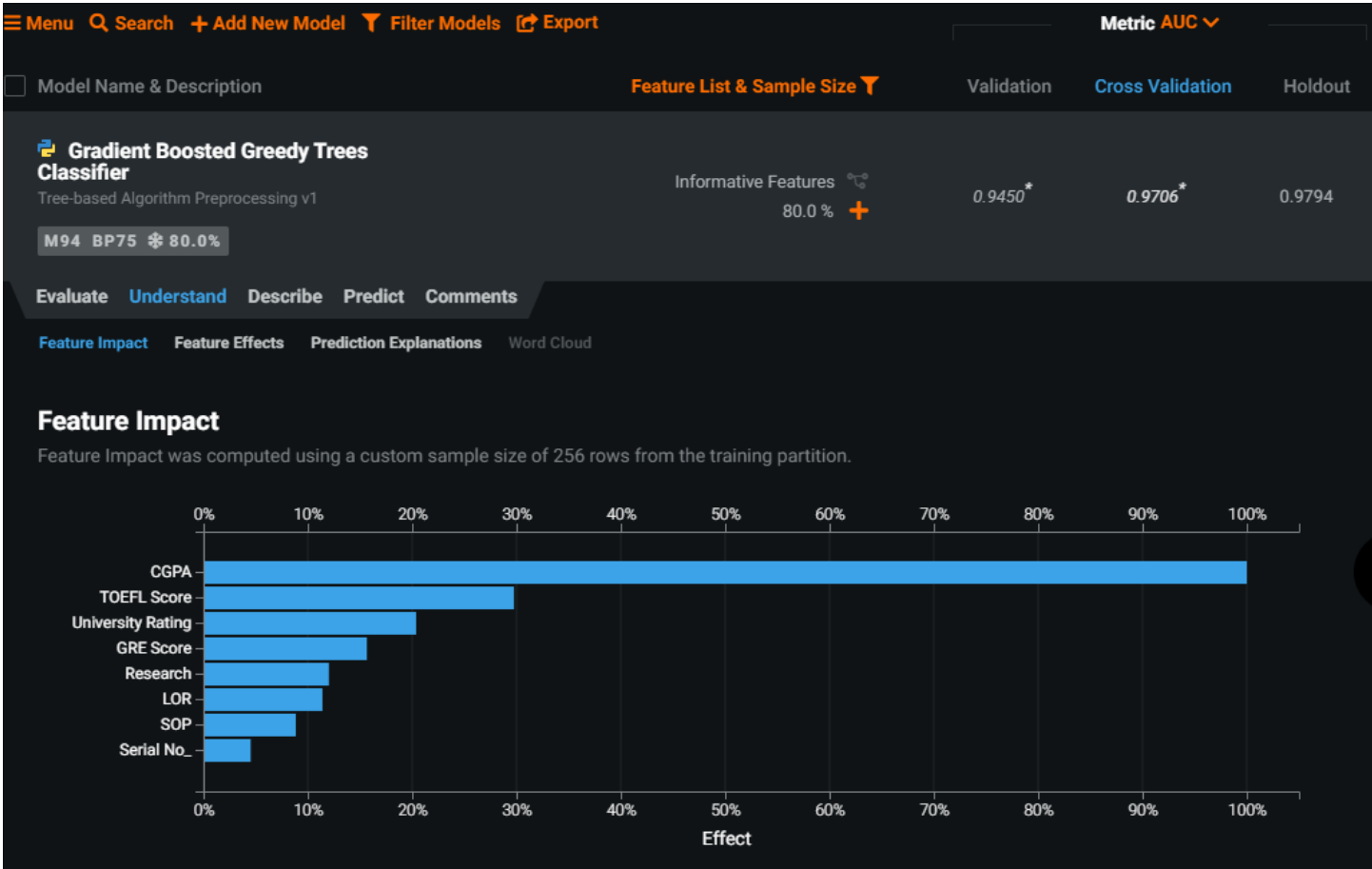| 📄 MODEL FILE SIZE | 🕐 PREDICTION TIME | ⚖ SAMPLE SIZE |
|---|---|---|
| **0.606 MB** | **6.4298s** ⚠ | **320 rows** |
| | Time to score 1,000 rows | Training 320 rows |
| | | Test 64 rows |

Here is the model evaluation results by using ROC Curve:



Here is the feature impact:

Next, I'll do the advanced tuning. The following screenshot shows the values of each parameter on the original Gradient Boosted Greedy Trees Classifier:

Model Name & Description

Advanced Tuning

**Parameters**

◉ New Search    ○ Searched    ○ Best of Searched

## Prediction Model Parameters

### Gradient Boosted Greedy Trees Classifier ⓘ

**learning_rate**

0.02

**max_depth**

None

**max_features**

0.1

**max_leaf_nodes**

3

**min_samples_leaf**

2

**min_samples_split**

5

**n_estimators**

500

**random_state**

1234

**subsample**

1.0

**Tuning parameters:**

I made some changes to max_depth, min_samle_leaf, and n_estimators.

Model Name & Description

**Advanced Tuning**

**Parameters**

⦿ New Search   ◯ Searched   ◯ Best of Searched

**Prediction Model Parameters**

**Gradient Boosted Greedy Trees Classifier** ⓘ

**learning_rate**

> 0.02

**max_depth**

> 16

**max_features**

> 0.1

**max_leaf_nodes**

> 3

**min_samples_leaf**

> 5

**min_samples_split**

> 5

**n_estimators**

> 200

**random_state**

> 1234

**subsample**

> 1.0

Finally, I got a better AUC scores on the training set (0.9716 vs 0.9706). For this model, the Undergraduate GPA is the most important feature. Next are TOEFL score and GRE score. The Statement of Purpose and Letter of Recommendation Strength is the least important feature in this model.