**Team 9**

George Mason University

**Analysis and Modeling of Inmate Recidivism Forecasting**

**Team Members:**

Kevin Uruena

Ryan Vavra

Lucien Wobembong

# Table of Contents

Contents

The Project Title: Analysis and Modeling of Inmate Recidivism Forecasting

Author's Name(s): Ryan Vavra, Kevin Uruena, Lucien Wobembong

Team #: 9

Header:
- **Course Professor's Name**: Dr. Eddy Zhang
- **Course Name, Number, and Section #**: AIT-580-001 (Summer 2021)
- **University Name**: George Mason University

# 1 Introduction

## 1.1 Problem

Recidivism "refers to a person's relapse into criminal behavior." There is a high rate of recidivism among released prisoners. From those released in 2005 (401,288), there was an average of 5 arrests per prisoner (1,994,000) within a 9-year period. Forty-four percent of those released in 2005 were arrested during the first year following release.

## 1.2 Objective

Team 9 wants to create a model that will help justice departments determine the likelihood of an inmate recidivating within 3 years, and if an inmate will recidivate at all.

# 2 Data Encoding and Handling of Null Values

Team 9 decided to conduct initiative data exploration by replacing null values with zeroes for drug related attributes. We decided on this approach since we wanted to assume the principle of 'innocent until proven otherwise.' Replacing with mean could have added values for some inmates who did not test positive. In addition, we performed label encoding for some of the attributes. This meant converting some of the categorical values into numeric values. For example, in the original dataset, gender was either M or F, but we converted that to 0 (M) and 1 (F). Race was converted from black or white to 0 for black and 1 for white. To make prediction easier, we decided to combine "Recidivism_Within_3years," "Recidivism_Arrest_Year1", "Recidivism_Arrest_Year2," and "Recidivism_Arrest_Year3." The "Encoded Dataset" file has the complete list of all the original and their converted attributes.

Encoded Dataset

# 3 Data Visualization

We utilized Tableau, a data visualization software to make our dataset interactive. We chose to make it interactive so audience members who are not familiar with data analysis terminologies, methodologies or do not care for models and predictions could still understand the data and our project. The simplicity of Tableau is depicted in figure 1 and 1.1 below.
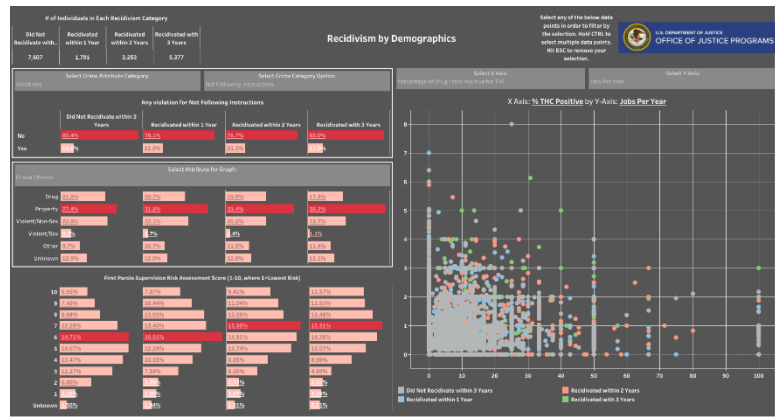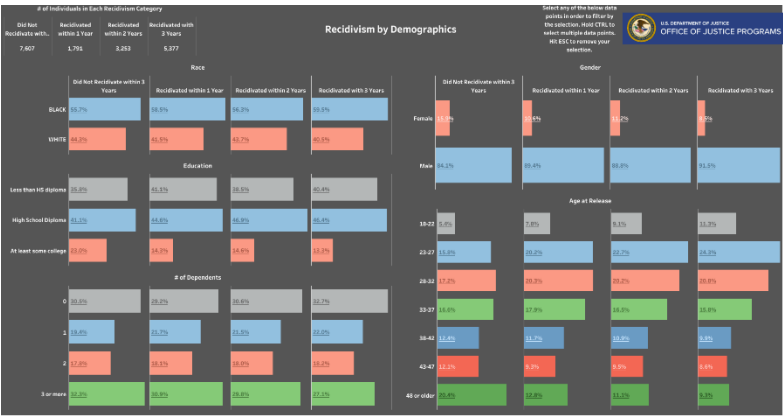


*Figure 1: Recidivism by Demographic*



*Figure 1.1: Recidivism by Demographic*

# 4 Improving Model Accuracy

## 4.1 Removing Attributes

We went further by looking at the most important attributes in the dataset by running a random forest, which is depicted in figure 1. At first, we considered dropping any features below 0.02 (which improved our model by 0.1%), but after experimenting with various models and algorithms, we decided keeping all attributes was best for higher accuracy model. We also decided to develop two different models. The first model was indented for multiple classification, separating the data into 4 different bins corresponding to no recidivism, recidivism within 1 year, 2 years, and 3 years. The second model was indented for binary classification, with the possible values being whether an individual will or will not recidivate.
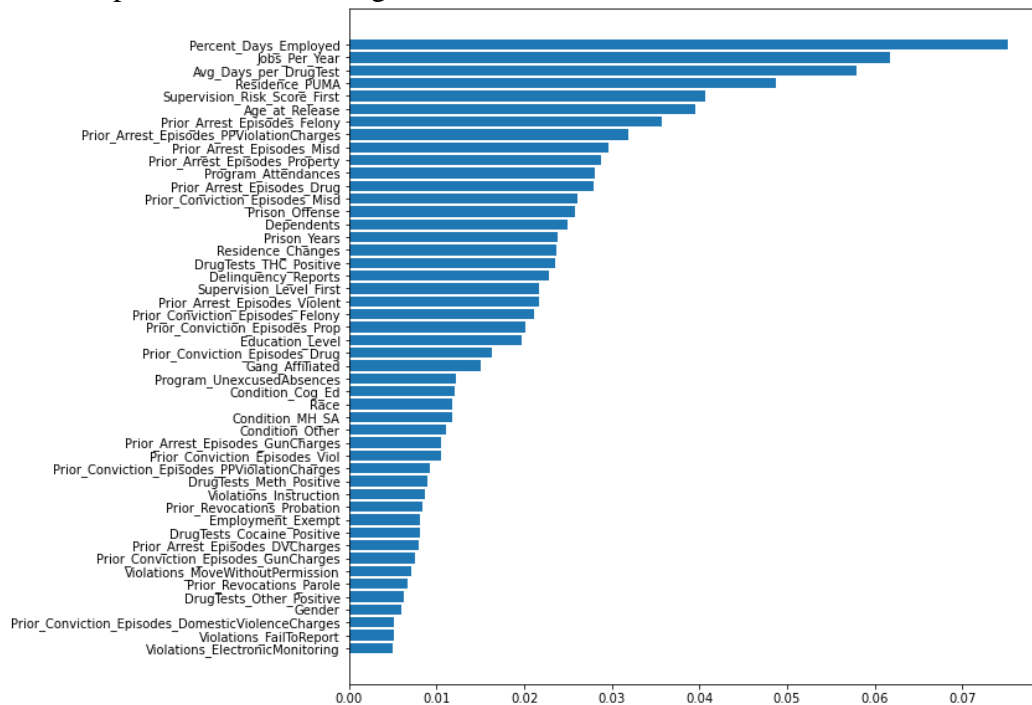


*Figure 2: Features Importance*

## 4.2 Sage Maker

Sage Maker is an AWS machine learning platform that allows the machine learning models to be easily developed and deployed. We decided to run autopilot mode in SageMaker studio to build and confirm the best performing algorithm. SageMaker selected XGBoost, trained, and tuned a model faster than our personal computers. We used the selected hyperparameters to run in local, confirming a similar performance. The model was saved in S3 and is ready for deployment in an EC2 instance.

## 4.3 Exhaustive Grid Search

The switch from attempting to classify four categories to trying to classify a boolean variable yielded an immediate jump in model accuracy. Despite the immediate improvement we were using only the default parameters of RF Classifier, besides n_estimators, so we hoped changing the parameters would yield an even higher accuracy score.

In addition to Sage Maker, we utilized the GridSearchCV function, imported from sklearn.model_selection, to tune our Random Forests Classifier Model. GridSearchCV is a function that allowed us to test a variety of hyperparameters and return us the set of hyperparameters that would yield the highest accuracy score.

Unfortunately, despite testing a wide variety of hyperparameter combinations, our model only improved 0.1% with the newly determined hyperparameters from the Grid Search.

# 5  Accuracy Results

To see which algorithms would be best for our dataset, Team 9 decided to try different algorithms and calculate the accuracy of each. As mentioned above, the 2 models developed were for binary, and multiple classification of inmate recidivism.

| Algorithms | Definition | Multiple Classification | Binary Classification |
|---|---|---|---|
| Random Forest | Uses averaging to improve accuracy and prevent over-fitting | 59.0% | 73.0% |
| Gradient Boosting | Converts weak learners into strong learners | 60.0% | 74.0% |
| XGBoost | Like Gradient Boosting, but improved for performance | 60.0% | 75.0% |

# 6  Keywords

The attached file has all the attributes and their description

NIJ Recidivism Data
Codebook.xlsx

# 7  References

https://www.simplilearn.com/10-algorithms-machine-learning-engineers-need-to-know-article

https://www.geeksforgeeks.org/ml-one-hot-encoding-of-datasets-in-python/

https://data.ojp.usdoj.gov/stories/s/daxx-hznc

https://towardsdatascience.com/understanding-gradient-boosting-machines-9be756fe76ab

https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

https://docs.aws.amazon.com/sagemaker/latest/dg/whatis.html