



# XGBoost Algorithm

Jared Mindel, Isaac Lee,  
Jeannine Hall, and Jed Dryer

# Process Overview

- Selected the dataset
- Prepared the data
- Built and tested the model/tuned hyperparameters
- Compared XGBoost to Logistical Regression
- Issues we encountered
  - Git control within group
  - Long runtimes and crashing computers

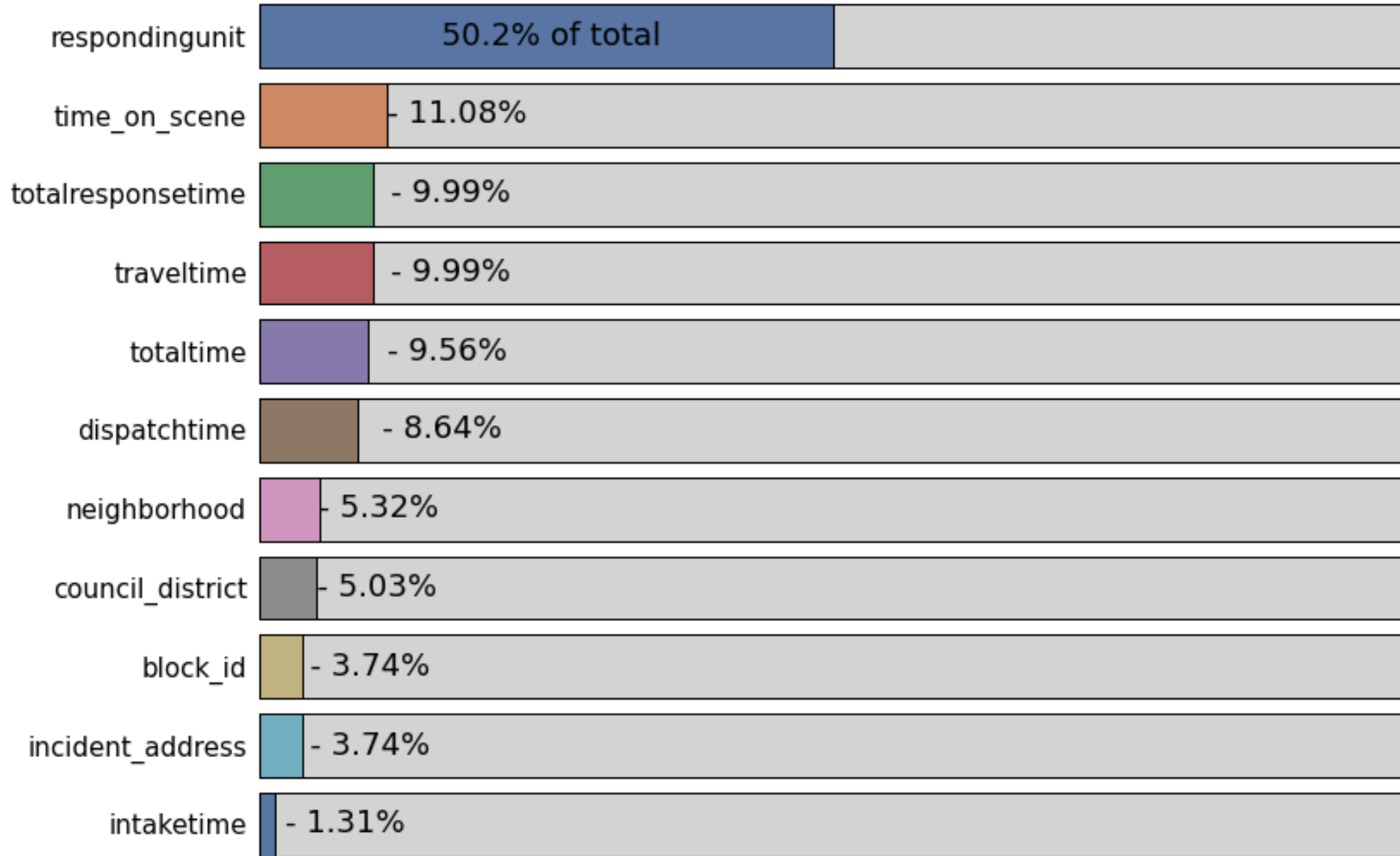
# Selecting the Data

- Datasets we considered:
  - Books
  - Diabetes
  - Cars 93
  - Airports
  - Spotify Songs
  - Meteorite Landings
- Our final choice: 911 Service Calls in Detroit

# Preparing the Data

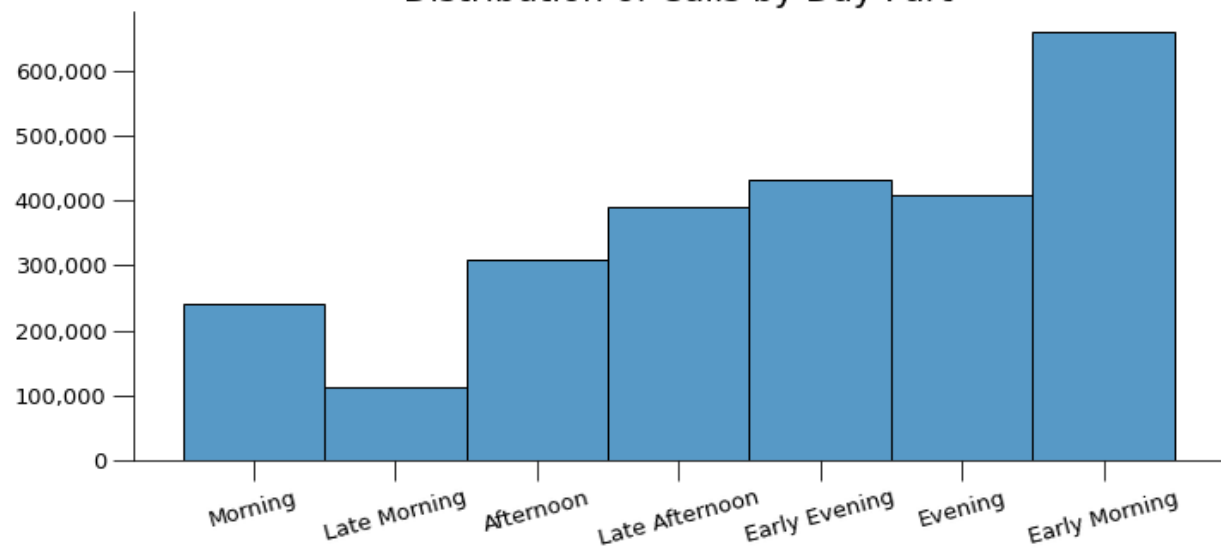
- Determined where the null values were
- Removed unnecessary columns
- Filled in missing values in Priorities column
- Created Weekday and Day\_Part columns
- Dummied categorical data
- Did not normalize data because XGBoost uses decision trees

## Features with Null Values

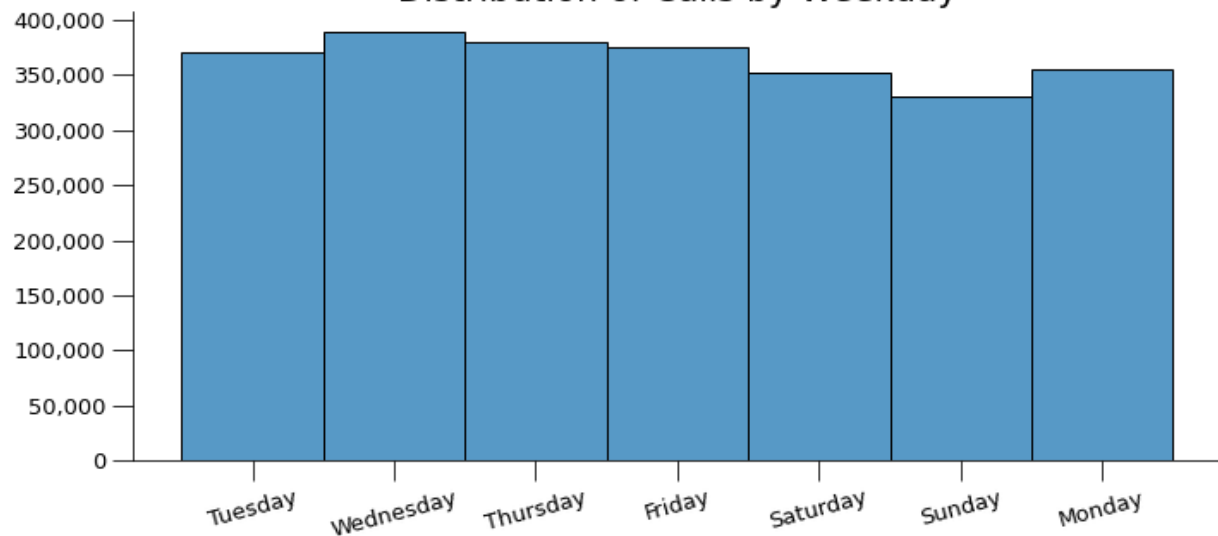


% of Values that are Null

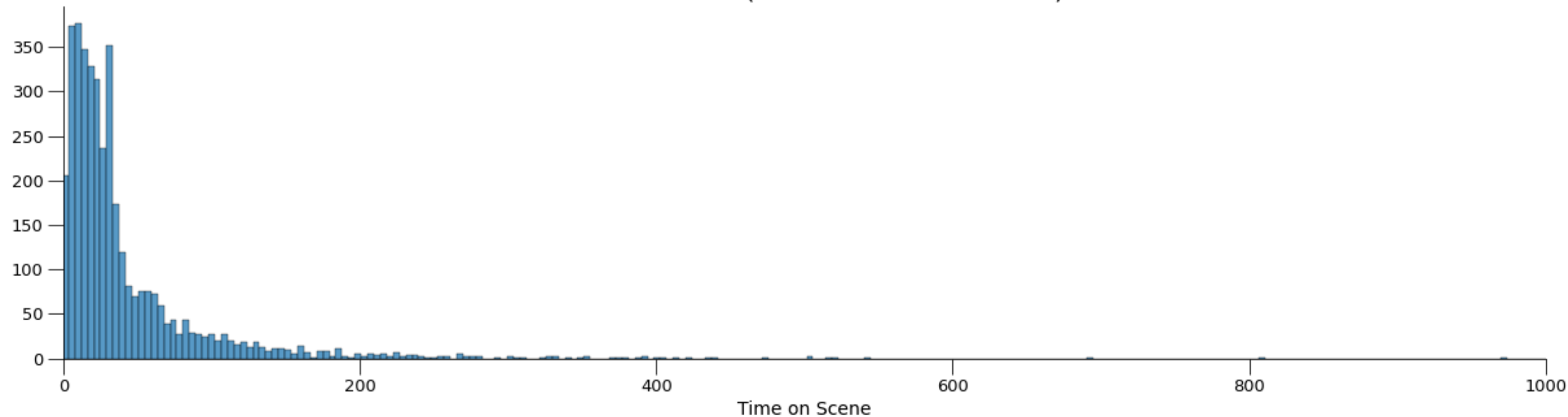
### Distribution of Calls by Day Part



### Distribution of Calls by Weekday



### Time on Scene (distribution of a subset)



# Running the model

Tried several different varieties of XGBoost

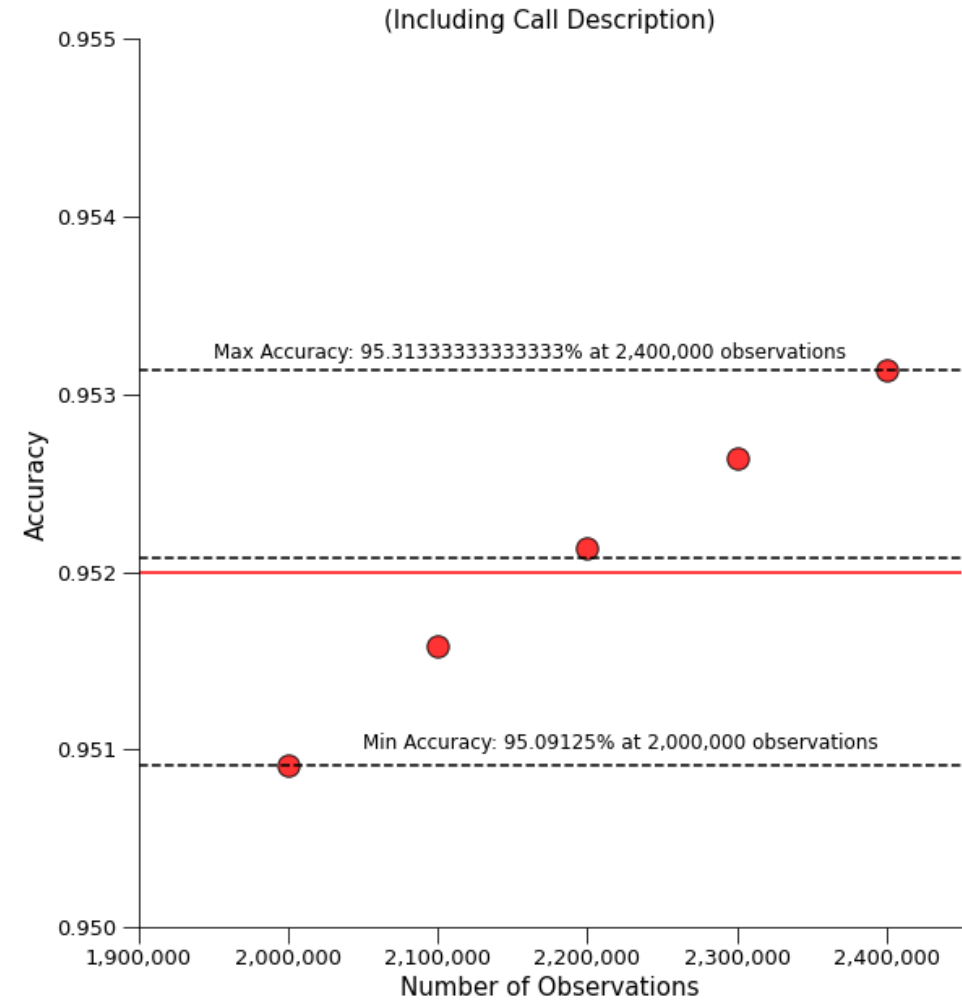
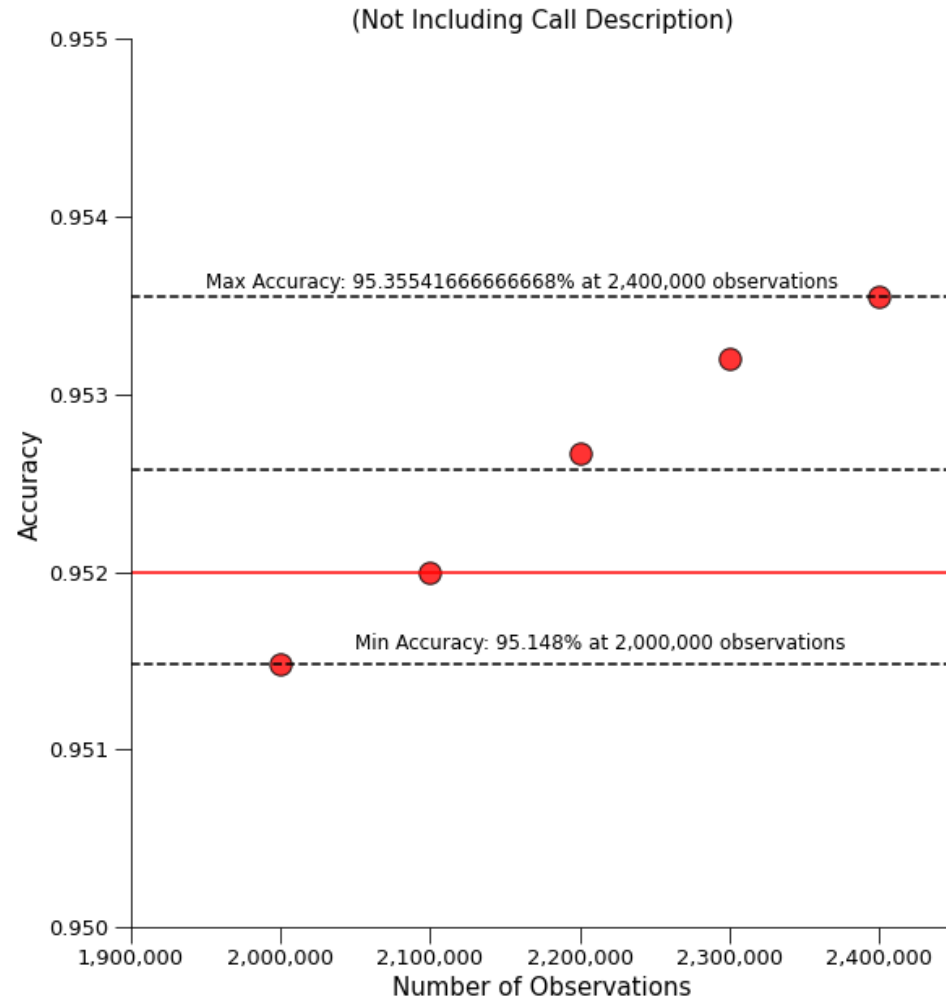
- DMatrix
- Classifier
- Regressor
- Booster

Settled on XGBoost Classifier

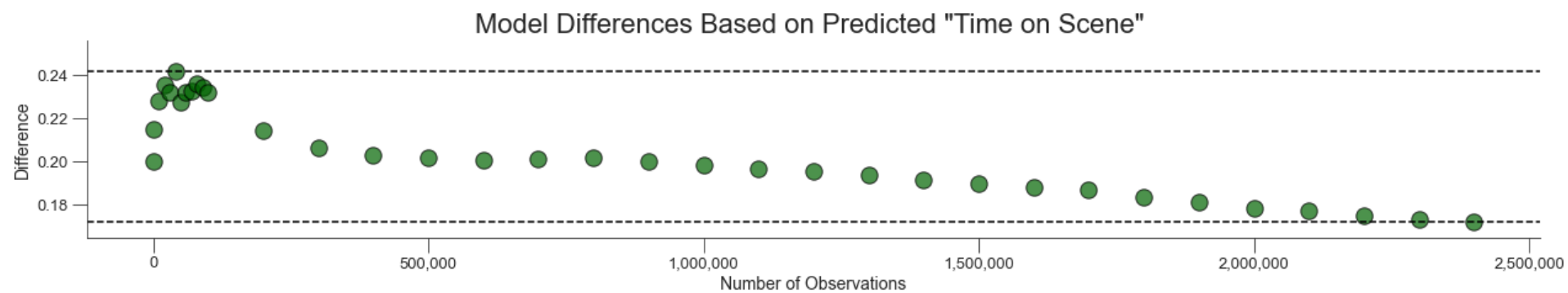
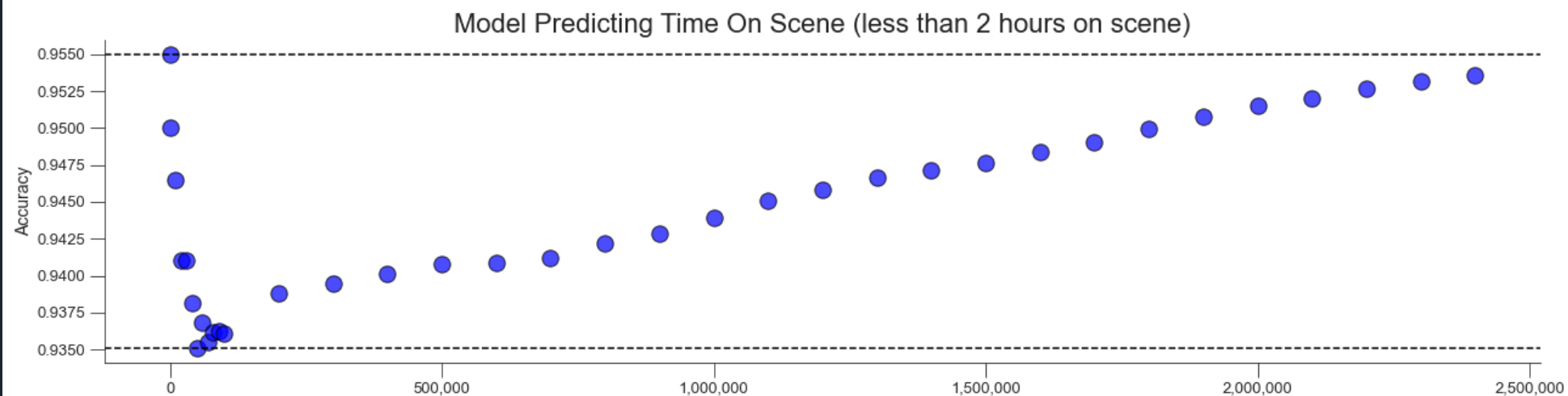
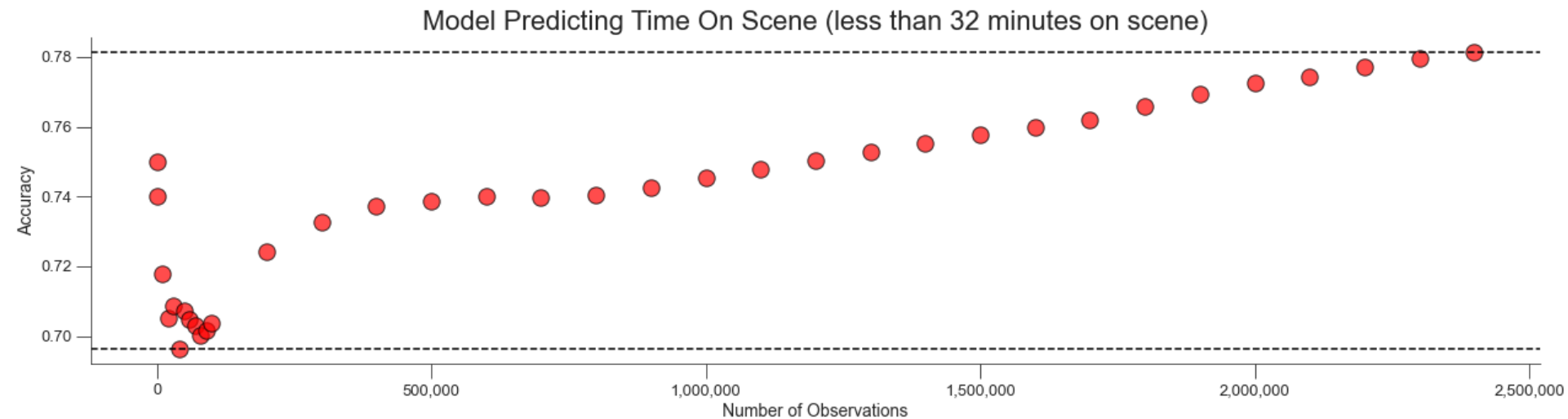
- Initially tested two-hour "time on scene"
- Changed our target to the average of the time on scene

# Running the model

Accuracy vs. Number of Observations



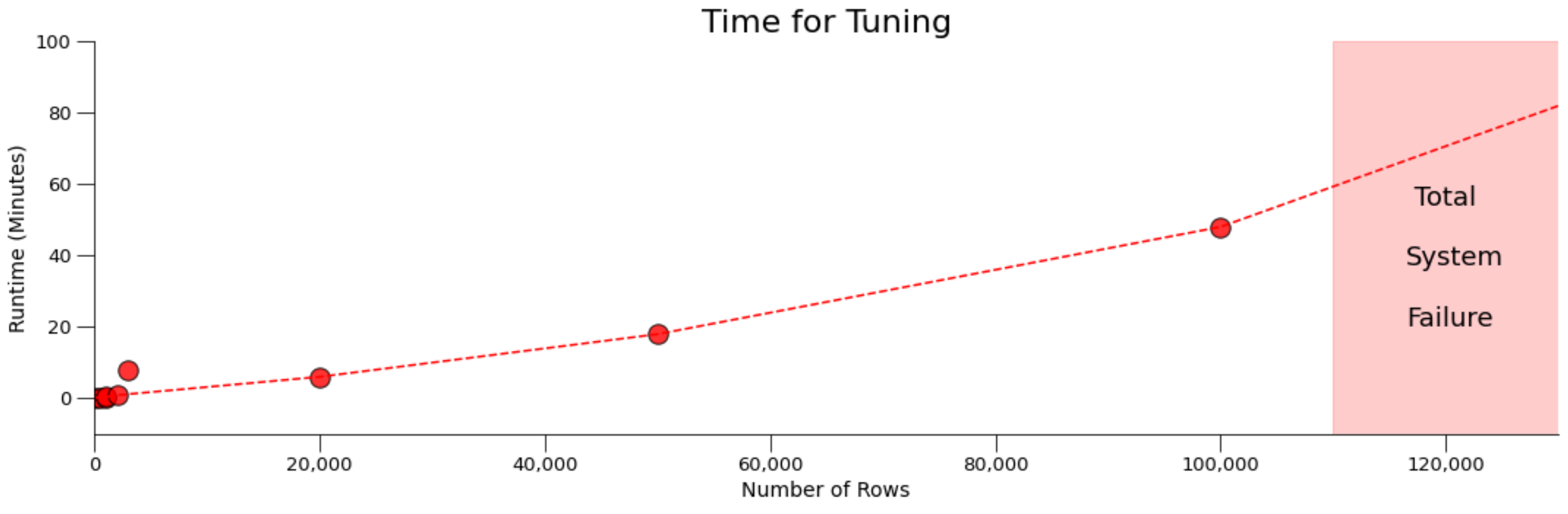




# Tuning the Model

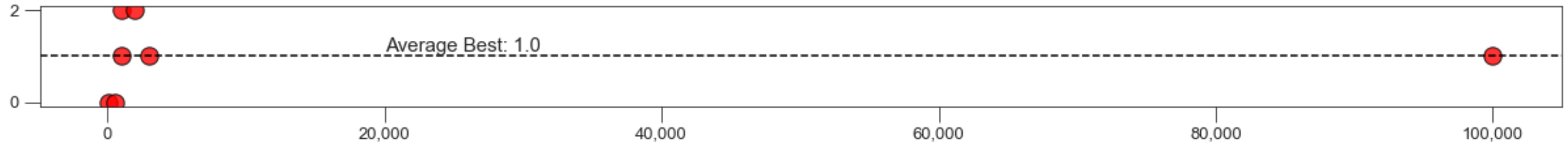
- Cross Validation with GridSearchCV
  - Passed in Hyper-Parameters and evaluated the results
  - System Crashes above 100,000 observations
  - During tuning, the best parameters from each round are preserved as the sole value for that parameter.
    - Used the average of the returned "best fit" in our final model.

# Tuning with GridSearchCV

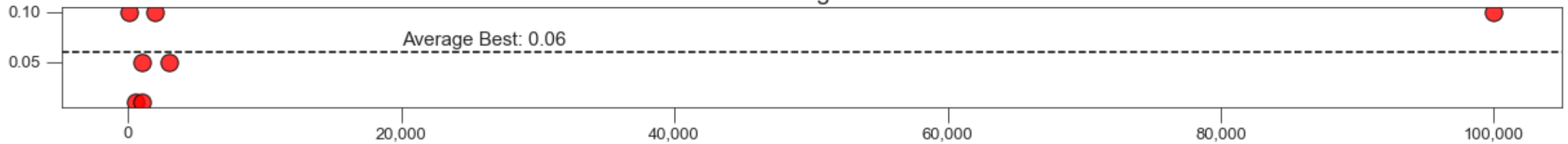


## Optimal Parameters (average of results)

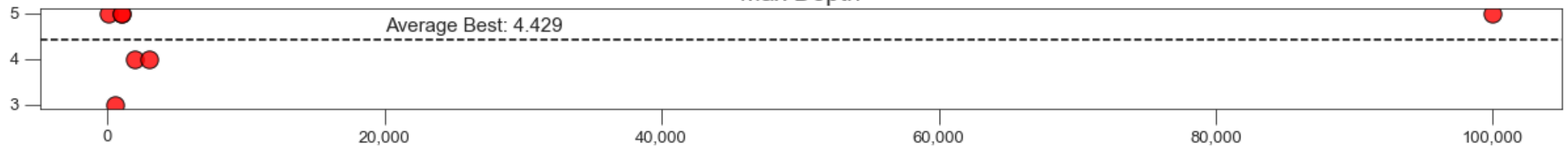
Gamma



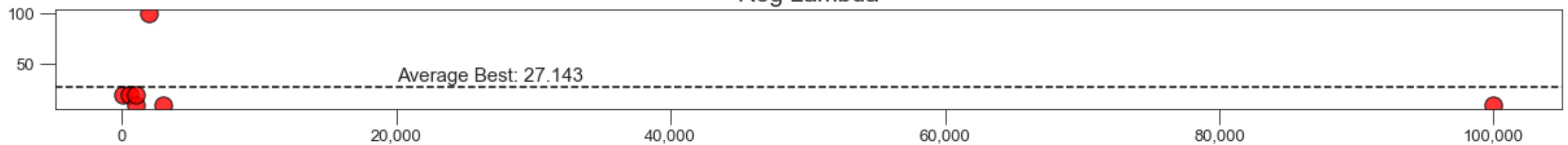
Learning Rate



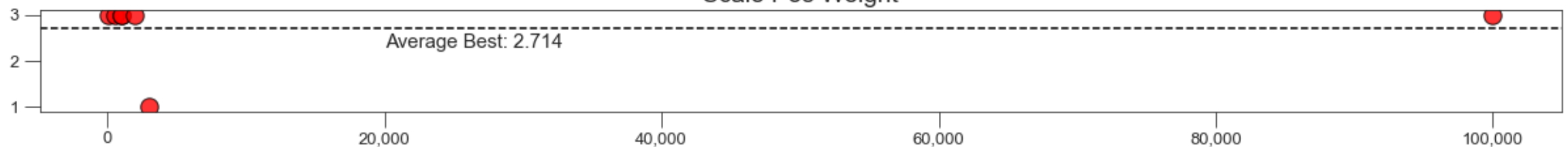
Max Depth



Reg Lambda



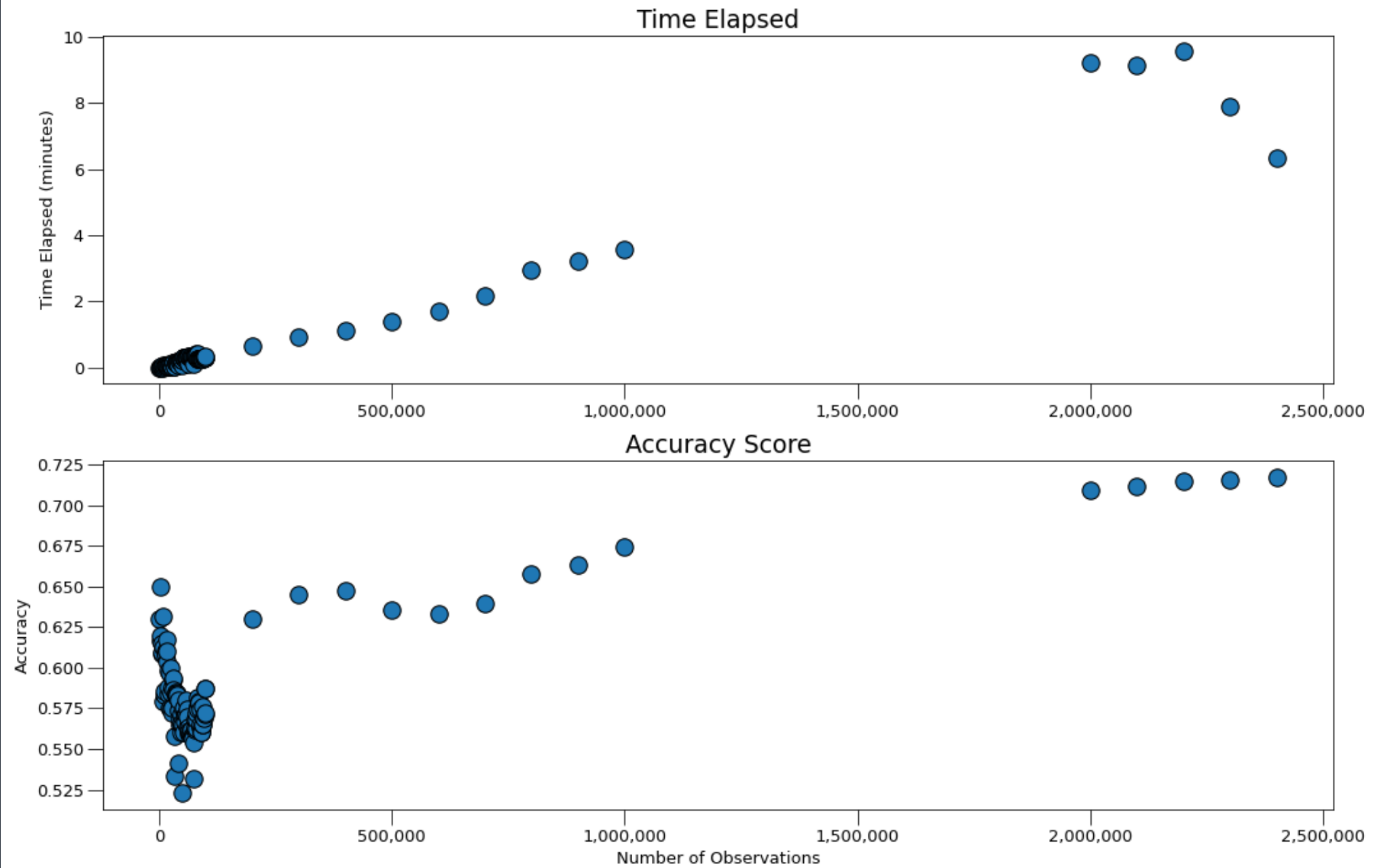
Scale Pos Weight



# Model Performance Improvements

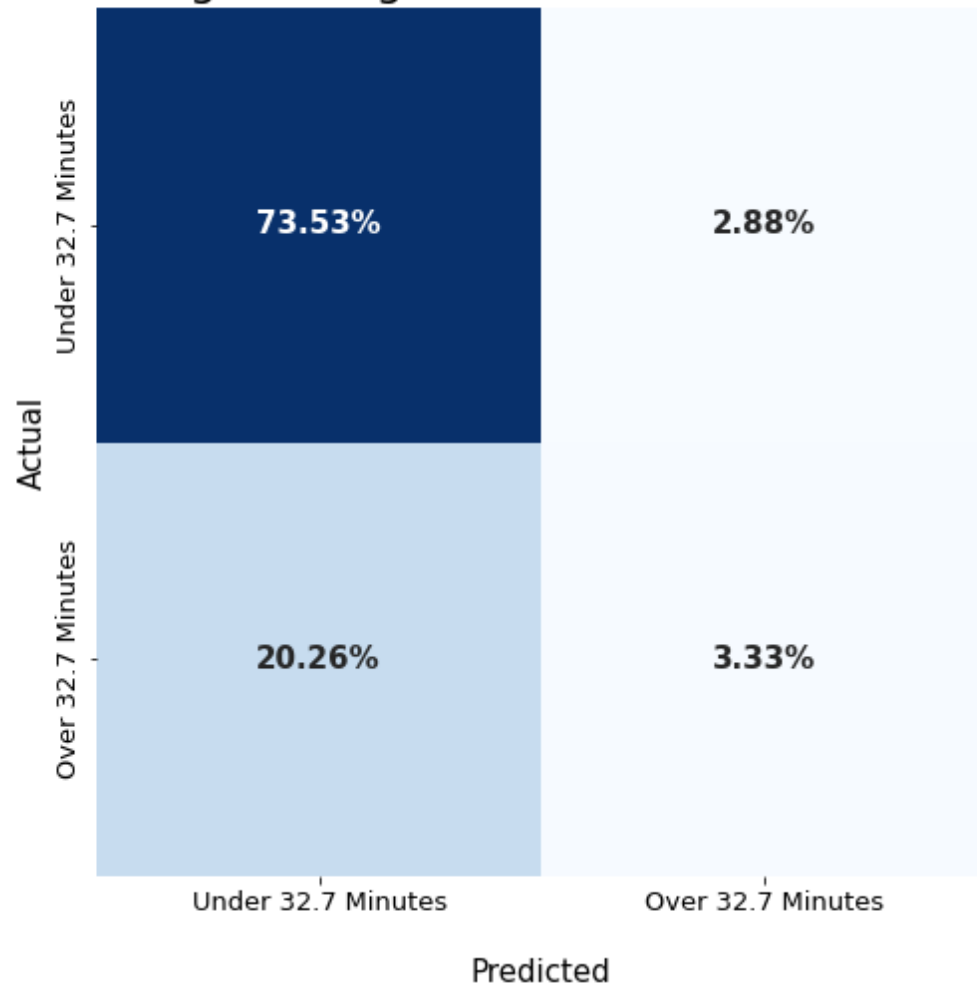
- Iterate fewer times – for sake of speed
- Boost memory performance
- Removed Call Description from the model
- Adjusting the random state may change the accuracy of the results
- Reverting to the pre-tuning hyper-parameters.

## Optimized Model Results

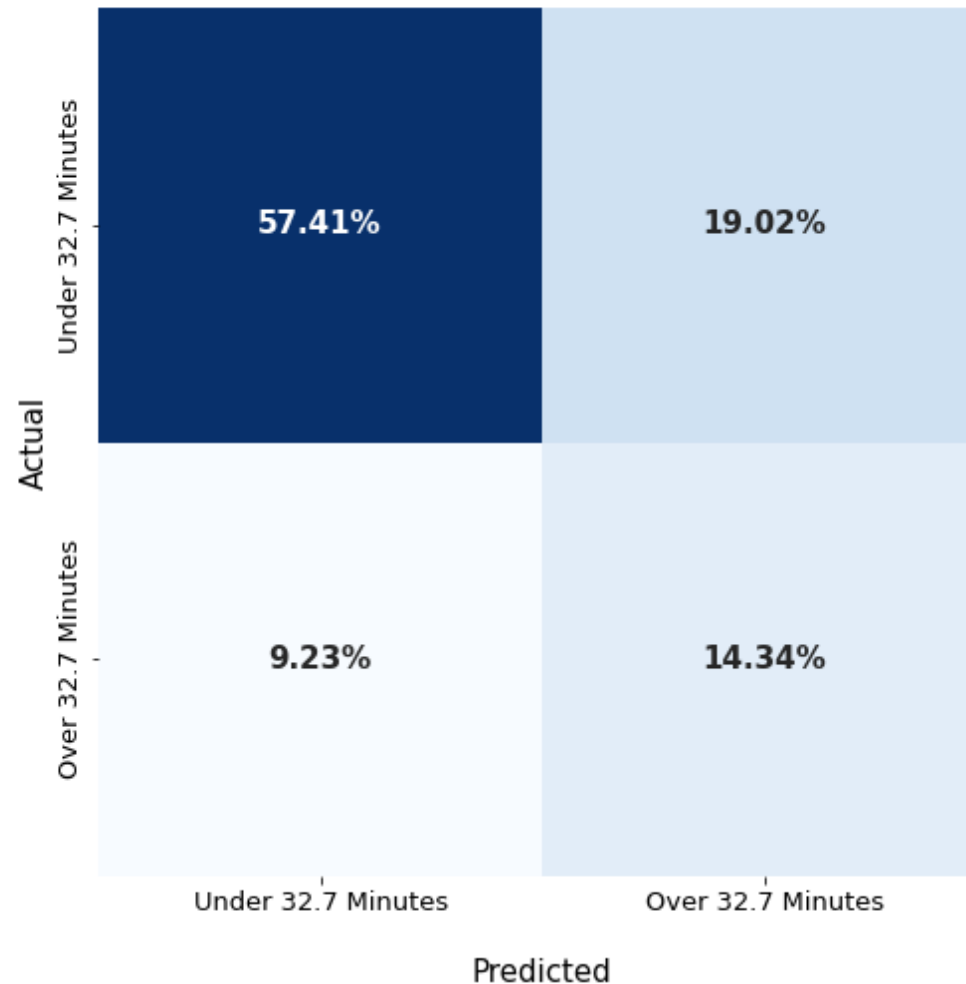


# Model Comparison: Confusion Matrices

Logistic Regression Confusion Matrix



XGBoost Confusion Matrix



# Questions?

