

Data Mining

CSCI-6674

**Analytics and Insights Team**

**Project Report Phase-5**

Master's in Computer science



University of New Haven

TAGLIATELA COLLEGE OF ENGINEERING, West Haven, CT

Submitted to:

Prof. Dr. Reza Sadeghi

Table of Contents

<b>Project Report: Phase – 5</b> .....	3
<b>Modeling Techniques</b> .....	5
<b>Model Parameters and Hyperparameters</b> .....	6
<b>Performance Metrics</b> .....	8
<b>Conclusion</b> .....	12
<b>References</b> .....	13

## Project Report: Phase – 5

### Team Details

Team Name : Analytics and Insights Team

### Team Members

1. Vyshampa Maringanti : [vmari2@unh.newhaven.edu](mailto:vmari2@unh.newhaven.edu) (Project Lead)
2. Akhila Gade : [agade1@unh.newhaven.edu](mailto:agade1@unh.newhaven.edu) (Team Member)
3. Bagyasree Chitikela : [bchit1@unh.newhaven.edu](mailto:bchit1@unh.newhaven.edu) (Team Member)

### Repository Address

<https://github.com/vmari2/Crime-Prediction-of-NYPD.git>

### Research Question:

Can we predict the occurrence of offenses for a given precinct within the next three months, based on historical crime data of the New York Police Department (NYPD) by using various classification techniques?

### Research Objective:

The main objective of this project is to create a predictive model using Negative Binomial Regression<sup>1</sup> to predict the number of crime occurrences/offences for a given precinct of the NYPD<sup>2</sup>. This prediction will aid the police department in determining the required patrolling manpower per precinct.

### Dataset Description:

The dataset covers the data of all crimes reported to the New York Police Department from 2006 to the end of last year i.e., 2019. This is a large dataset consisting of 6.98 million

records with 35 attributes. Given, the large size, we could filter data to a more manageable size and per our project requirement.

**Data set Source:**

<https://www1.nyc.gov/site/nypd/stats/crime-statistics/citywide-crime-stats.page>.

## Modeling Techniques

Modeling is of various kinds. Predictive Modeling is used to analyze the past data and predict the future outcome. In our project we implemented modeling on NYPD Dataset. NYPD maintains monthly crime data based on police incident reports across New York. The data is saved on the NYPD website with each month corresponding to an Excel file. Different modeling algorithms in Python have been used for the prediction.

There are various modeling techniques which can be used, some of the modeling techniques we used are as follows:

Modeling Technique	Accuracy	Precision	Recall
Logistic Regression	99.96	1.0	1.0
Random Forest	100	1.0	1.0
XG Boost	100	1.0	1.0
K-Nearest Neighbor(KNN)	99.97	1.0	1.0
Decision Trees	100	1.0	1.0

### Hardware :

- **Processor:** Intel Core i5 – 8265U @ 1.60GHz 1.80GHz
- **RAM:** 8.00 GB
- **System Type:** 64-bit Operating System, x64-based processor

## Model Parameters and Hyperparameters

### Logistic Regression:

Logistic Regression has different parameters, but we used specific parameters such as ‘C’ and ‘penalty’. Penalty is used to specify the regularization. Usually ‘l2’ is the default penalty, which is used by most of the solvers, if we are not specifying any penalty then no regularization is applied. C parameter is used to control the penalty strength.

### Random Forest:

The parameters used in Random Forest are ‘n\_estimators’, ‘min\_samples\_split’, ‘min\_samples\_leaf’. The number of random features to sample at each split point is the most significant parameter. A log scale of 10 to 1,000 might be good value. Until no further change is seen in the model, the number of trees should be increased.

### XG Boost:

In XG Boost, the parameters used for the model are ‘n\_estimators’, ‘max\_depth’, ‘reg\_lambda’ and ‘learning\_rate’. The **n\_estimators**: Number of gradient boosted trees. Equivalent to number of boosting rounds, **max\_depth**: Maximum tree depth for base learners, **learning\_rate**: Boosting learning rate and **reg\_lambda** : L2 regularization term on weights.

### KNN:

The number of neighbors (n neighbors) is the most significant KNN hyperparameter. It measures values between 1 and 21. The optimal number of neighbors for our project are 19. We also have some parameters in KNN such as weights and metric.

**Decision Trees:**

The parameters used in Decision Tree are like Random Forest, but we have an extra parameter namely 'max\_depth' . If no max\_depth values are specified, then the nodes are expanded until all leaves are either pure or till they are less than the min\_samples\_split values.

## Performance Metrics

**Accuracy:** We got an accuracy of 99% for all our models for both training and testing data by using all features of the dataset. Below table represents the training and testing data accuracy of different models.

Modeling Technique	Train data accuracy	Test data accuracy
Logistic Regression	99.96	100
Random Forest	100	99.99
XG Boost	99.97	100
K-Nearest Neighbor(KNN)	99.74	99.99
Decision Trees	100	100

**Precision and Recall :** Precision and Recall are the evaluation metrics of a model. Precision helps us to know what proportion of the positive identifications are correct and Recall helps us to know what proportion of actual positives are correctly identified.

For the dataset and features we used, all the models had a precision and recall of 1.

### Feature selection and Feature importance graphs for each model :

**1. Logistic Regression :** The below are the F-score for each feature used in the model.

We can see that the highest score is recorded for Feature 2.

Feature: 0, Score: 0.16320

Feature: 1, Score: -0.64301

Feature: 2, Score: 0.48497

Feature: 3, Score: -0.46190

Feature: 4, Score: 0.18432

Feature: 5, Score: -0.11978

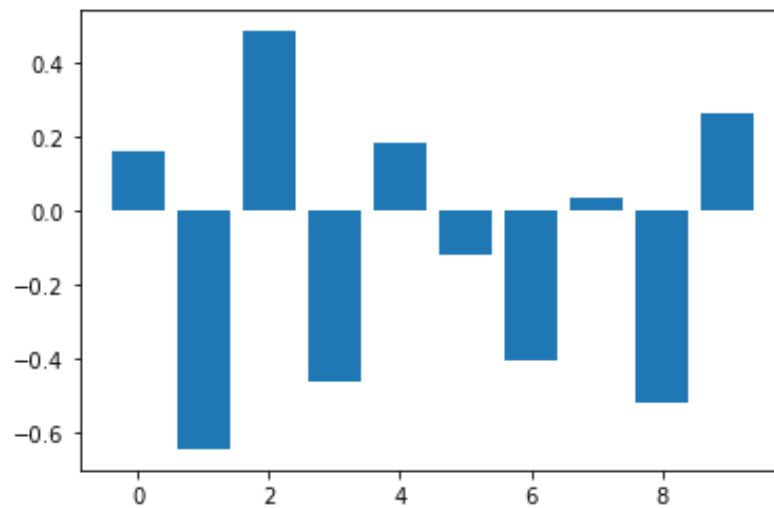


Feature: 6, Score: -0.40602

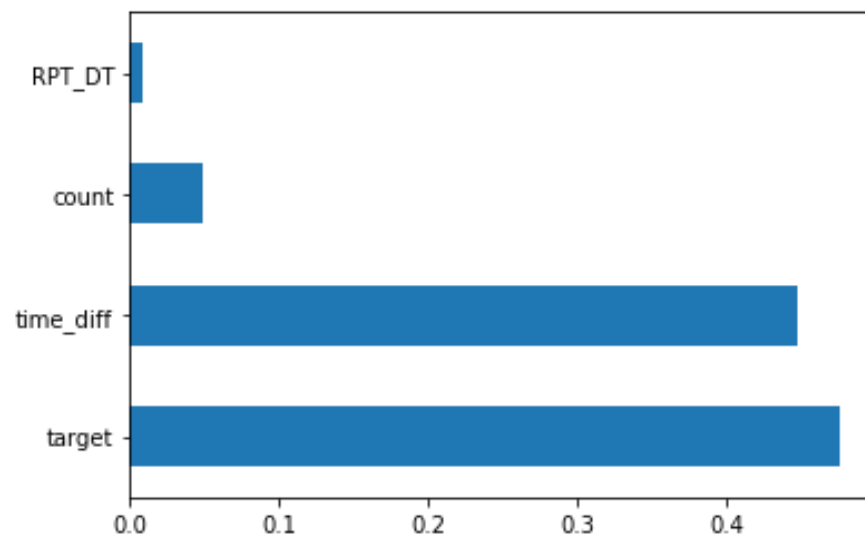
Feature: 7, Score: 0.03772

Feature: 8, Score: -0.51785

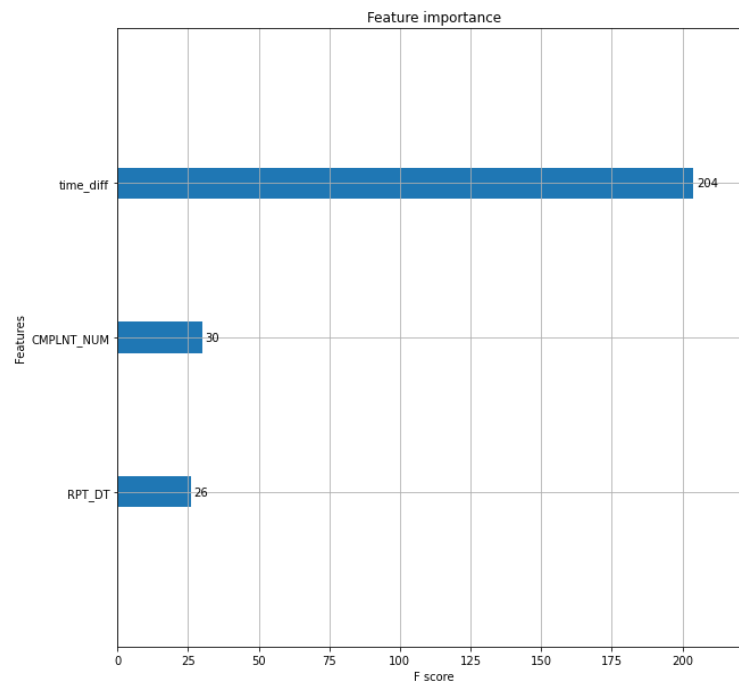
Feature: 9, Score: 0.26540



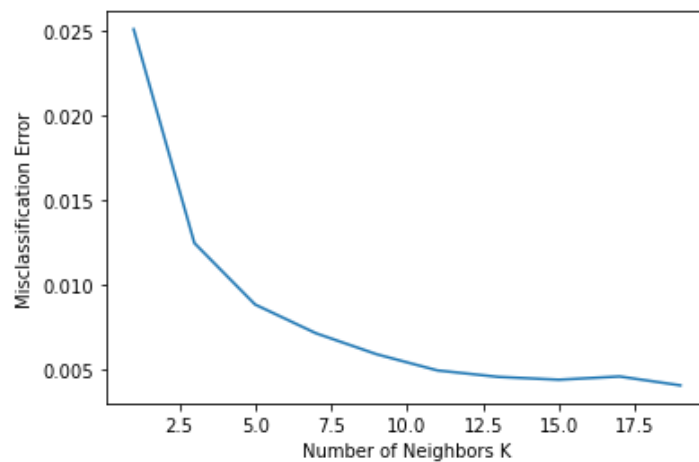
- 2. Random Forest :** In the below graph, we can see that the target has the highest score of more than 0.4



3. **XG Boost:** From the below graph it can be noted that the time\_diff has more importance than the other features.



4. **KNN :** In the below graph, we can observe that the misclassification error decreases as the number of neighbors k value increases.



**5. Decision Trees :** From the below values we can see that the highest score is recorded for the Feature 0.

Feature: 0, Score: 0.45619

Feature: 1, Score: 0.04764

Feature: 2, Score: 0.08236

Feature: 3, Score: 0.00109

Feature: 4, Score: 0.39992

Feature: 5, Score: 0.00118

Feature: 6, Score: 0.00121

Feature: 7, Score: 0.00145

Feature: 8, Score: 0.00155

Feature: 9, Score: 0.00112

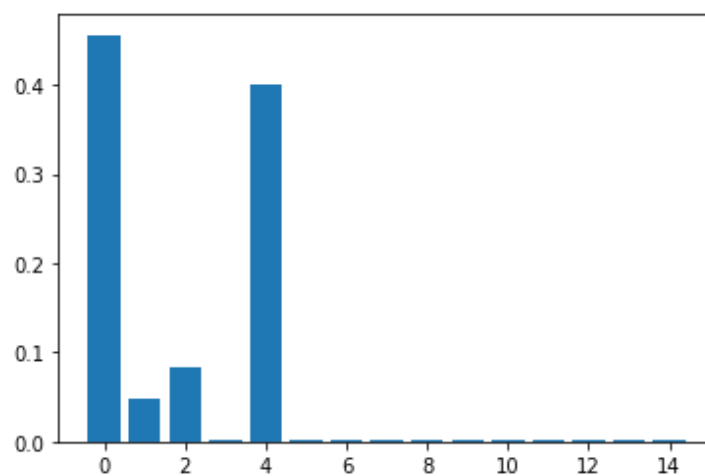
Feature: 10, Score: 0.00133

Feature: 11, Score: 0.00161

Feature: 12, Score: 0.00110

Feature: 13, Score: 0.00125

Feature: 14, Score: 0.00101



## Conclusion

The dataset used for the modeling includes 60,000 records with 24 features. The dataset is divided into training dataset (70%) and testing dataset (30%). The training dataset was used for training the models and the test dataset was used to test the accuracy of the model. By training our models with the training dataset we observed that all our models provide the accuracy of 99%. From this we can conclude that there is a higher possibility of crime in New York City.

## References

1. <https://machinelearningmastery.com/hyperparameters-for-classification-machine-learning-algorithms/>
2. <https://scikit-learn.org/>
3. [https://xgboost.readthedocs.io/en/latest/python/python\\_api.html#module-xgboost.sklearn](https://xgboost.readthedocs.io/en/latest/python/python_api.html#module-xgboost.sklearn)