



The University of Texas at Austin



# Chicago Crime Incident Analysis

## Project Report

INF 397 – Statistical Analysis and Learning w/ Prof. Varun Rai  
Spring 2018

### **TEAM MEMBERS:**

Milind Siddhanti (mss4376)  
Nimish Kate (nk8648)  
Prachi Singh (ps28755)  
Sanchit Singhal (ss84657)  
Shreshtha Shukla (ss83452)

# Introduction

As the mankind takes its steps towards digitization, the world eagerly follows. Information grows to become the most sought out resource. The ever-increasing presence of information can then be used to perform analysis and better our understanding of the world. In this constant race towards a better future, security is a important feature that presents an increasing concern. Cities are growing at a rapid rate, and so is the population. All these factors have led to an increase in crime, making it a priority for the authorities to identify crime-prone areas and take precautionary measures. To understand neighborhoods that need extra attention, it is imperative to identify zones with higher crime rates, predict recurring crimes, their frequencies, and analyze the time frames in which most crimes occur. A detailed analysis of the crimes would help the authorities understand the type of security measures that need to be put in place.

We attempt to explore the data to first find the most common types of crime, and sort them according to the neighborhoods. We attempt to understand the frequency of arrests for different crimes based on various features. By classifying crimes into whether or not arrests will occur, this project aims to assist law enforcement by providing a better understanding of arrest patterns in the city of Chicago. We will then analyze these findings and provide recommendations to local police authorities to improve policing practices and training.

## Problem Statement

This report will provide a statistical analysis of criminal activities in Chicago and attempt to classify arrests based on features relating to the crime's features.

## Data & Preprocessing

The dataset for the analysis is borrowed from kaggle.com. It reflects reported incidents of crime (with the exception of murders where data exists for each victim) that occurred in the City of Chicago from 2001 to 2017. Data is originally extracted from the Chicago Police Department's CLEAR (Citizen Law Enforcement Analysis and Reporting) system. In order to protect the privacy of crime victims, addresses are shown at the block level only and specific locations are not identified.

The dataset is a total of 1.85 GB. The entire data is split into 4 groups of specific time spans consisting of over 6,000,000 records.

- 2001-2004 - 1923517 records
- 2005-2007 - 1872346 records
- 2008-2011 - 2688712 records
- 2012-2017 - 1456715 records

The available data had inconsistencies in records. There were missing values for certain features, and erroneous values for others. A few observations did not follow a consistent format of values. The data was cleansed to eradicate these issues. Data generation was also performed for missing data (records with missing geospatial information but a known location were transformed to generate the missing entries). Many records had spaces, quotations and other punctuation, which made organizing the data difficult. These contraptions were removed to make the data more robust.

Observations from the first file were removed from the refined dataset to reduce computational resource requirements - due to slow model training, a design decision was made to only analyze data from 2005 onwards. Also, it should be noted that the year 2017, only had records for a few months and therefore significantly different aggregation numbers compared to the rest of the years.

The data has 22 columns, which will act as variables while building the models; they are listed below:

### Features

ID - Case ID (Unique Identifier)  
Case Number - Case number of the crime reported  
Date - Date of the occurrence of the crime  
Block - Area (block) of Chicago where the incident took place  
IUCR - Illinois union crime reporting code  
Primary Type - Top most classification of the crime  
Description - Description of the crime  
Location Description - Description of the location  
Arrest - Arrested or not (True or False)  
Domestic - Type of crime  
Beat - Police jurisdiction  
District - Police Jurisdiction  
Ward - Jurisdiction of the crime scene  
Community Area - Jurisdiction of the crime scene  
FBI Code - FBI code defining crime category  
X Coordinate - X coordinates of the crime scene location  
Y Coordinate - Y coordinates of the crime scene location  
Year - Year the case was filed  
Updated On - Date/Time of the update on the crime case  
Latitude - Location of the crime scene (Latitude)  
Longitude - Location of the crime scene (Longitude)  
Location - Location of the crime scene (Latitude+Longitude)

## Data Exploration

An exploratory analysis of the cleansed data revealed the following:

- The number of crime incidents encountered by the Chicago crime department per year.

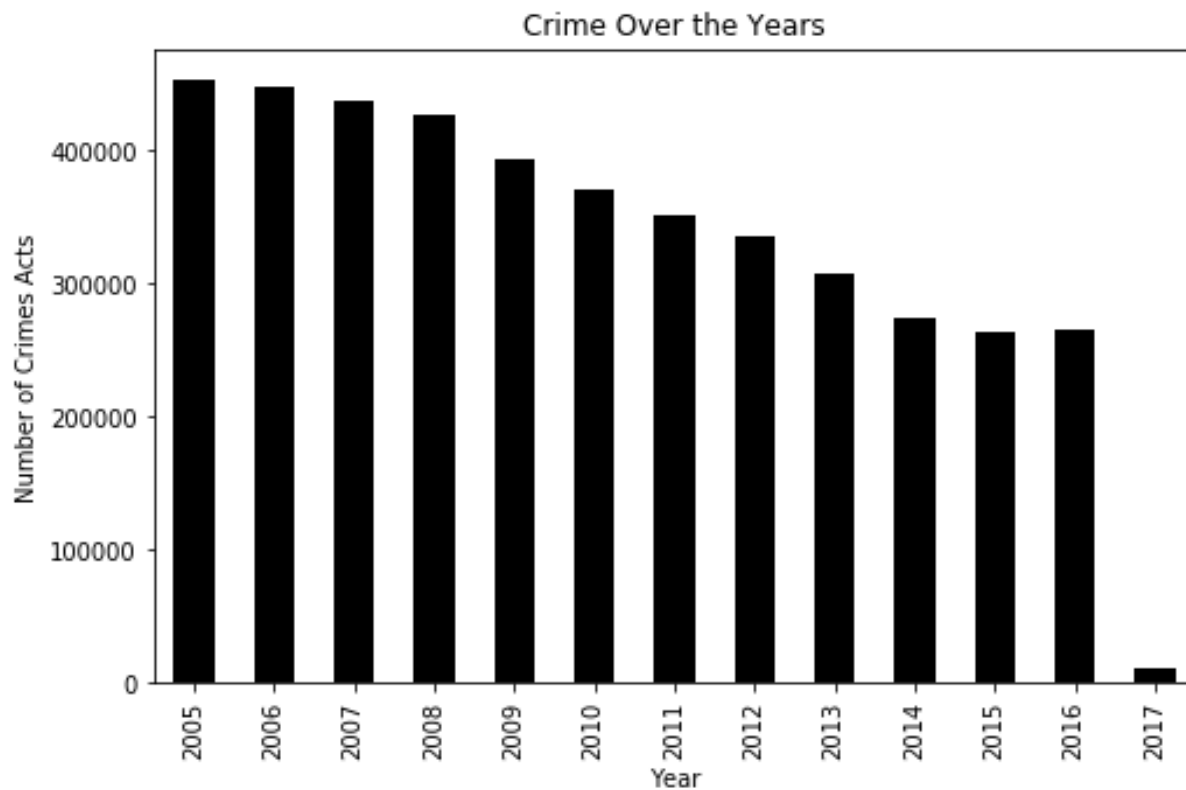


Fig. 1. Number of crimes reported each year

In Fig. 1, we can observe that the number of crimes reported in Chicago reduced gradually from 2005 to 2017. There was a significant drop in the crimes reported in the year 2009 and was almost constant in 2015 and 2016. As we do not have complete data for 2017, we cannot compare that with the data collected over the years.

- The relative number of crime incidents encountered by the Chicago crime department per month over the span of 13 years.

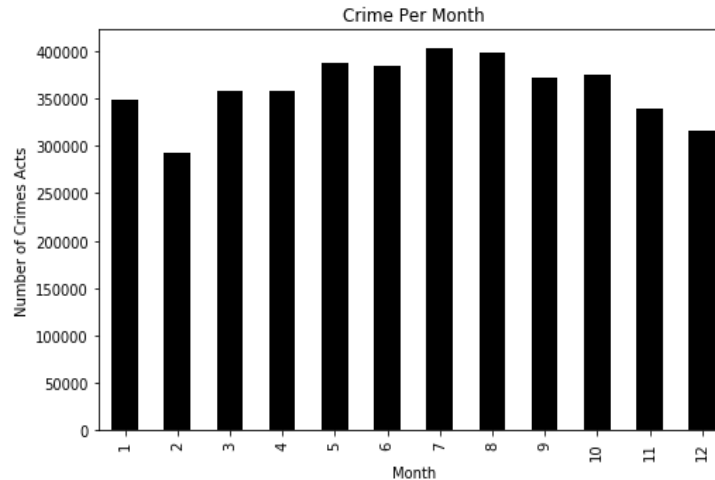


Fig. 2. Number of crimes reported per month

In Fig. 2, we can observe that the highest number of crimes were reported in July and the lowest in the month of February.

- The relative number of crime incidents encountered by the Chicago crime department per day of each month over the span of 13 years.



Fig. 3. Relative number of crimes reported every day per month

In Fig. 3, we can observe that the highest number of crimes were reported on the first of every month. There seems to be a drop in reported crime in the last three days of the month but the other days seem to have approximately the same number of reported crimes.

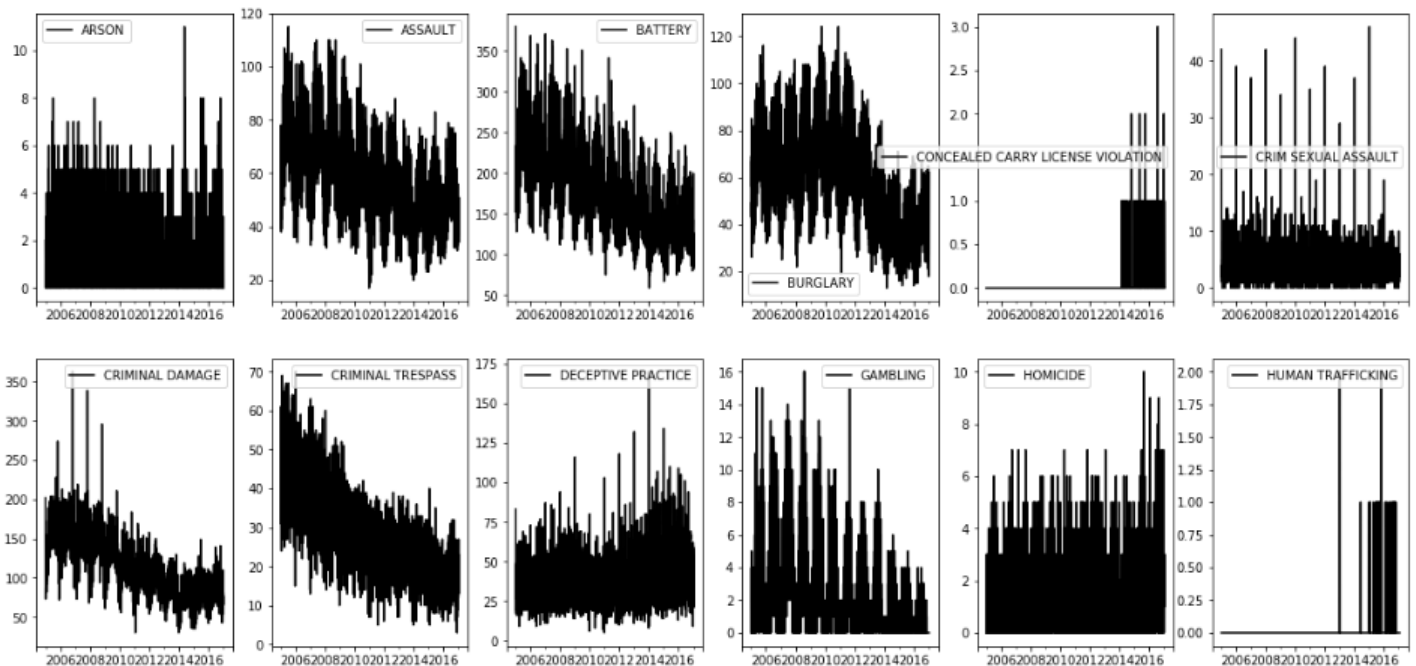
- The relative number of crime incidents encountered by the Chicago crime department every hour of the day over the span of 13 years.



Fig. 4 Number of crimes reported every hour per day

In Fig. 4, we can observe that the highest number of crimes were reported around 12 PM every day. The crimes were significantly lower during the early hours of the day and increased during the evening and night time.

- Types of crime per month



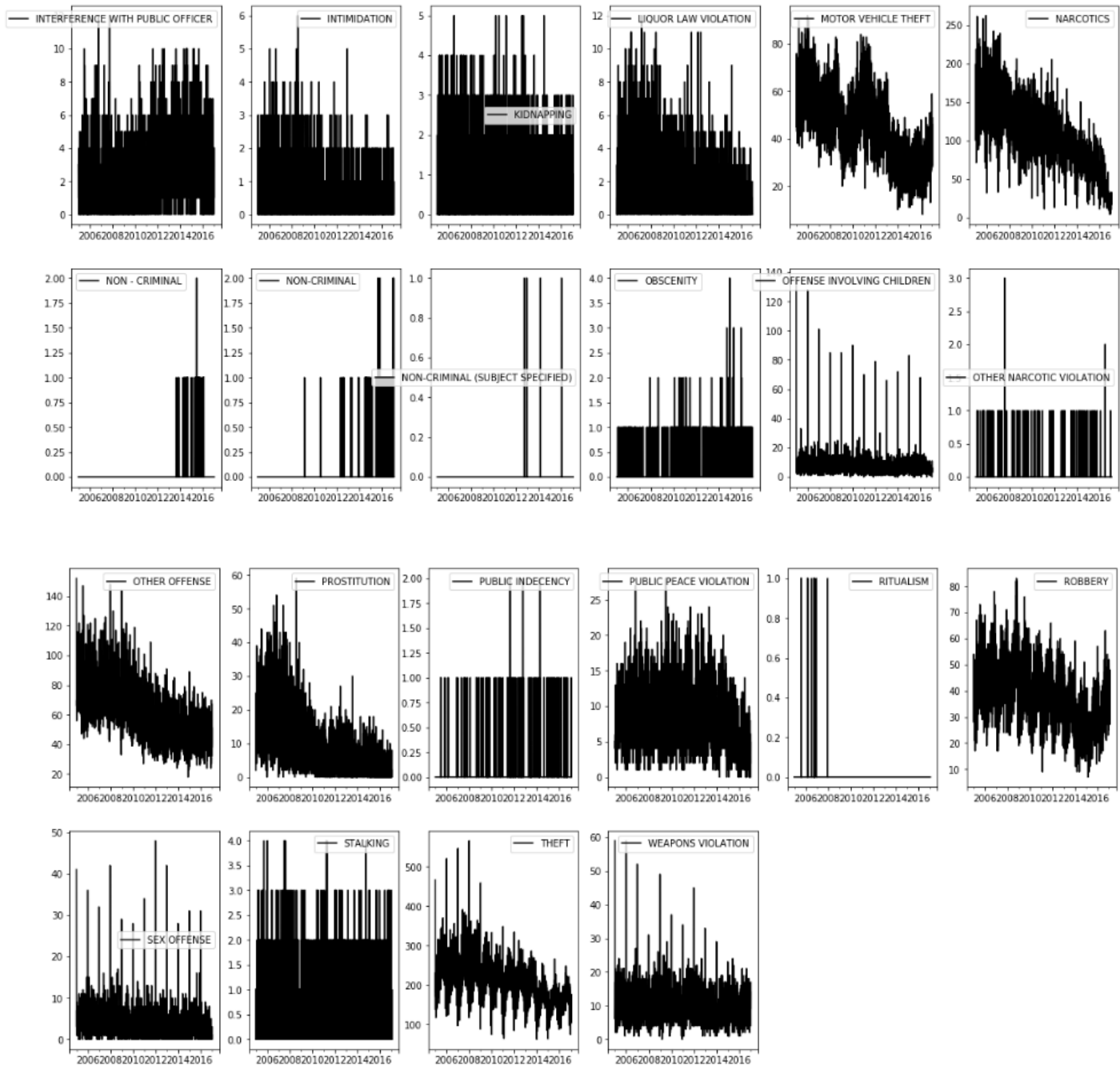


Fig. 5 Types of crimes reported per month

The number of crime incidents for the year 2017 are significantly less than those reported for the other years. A manual scan of the data revealed that the dataset only has records for January of 2017. We decided to exclude the data for the year 2017 noting its incomplete and incoherent nature. As it is shown in fig. 1, the crime reported was substantially higher in the year of 2008, and has steadily declined over the years. An exploratory analysis of this revealed that many crimes were reported at midnight on the first day of the year. This data appears to be unprecedented, and a possible explanation might be the way the crime department operates on the first day of the year. In order to reduce inaccuracies in the data, we shifted our analysis to focus on the years, rather than the time of a particular day.

Despite having a significantly large number of crime incidents from 2006 to 2010, the number of arrests did not increase to match the crime reports. The difference between the crime reports and arrests seems to perpetuate over the span of 16 years.

- The number of crime incidents reported for each type of crime.

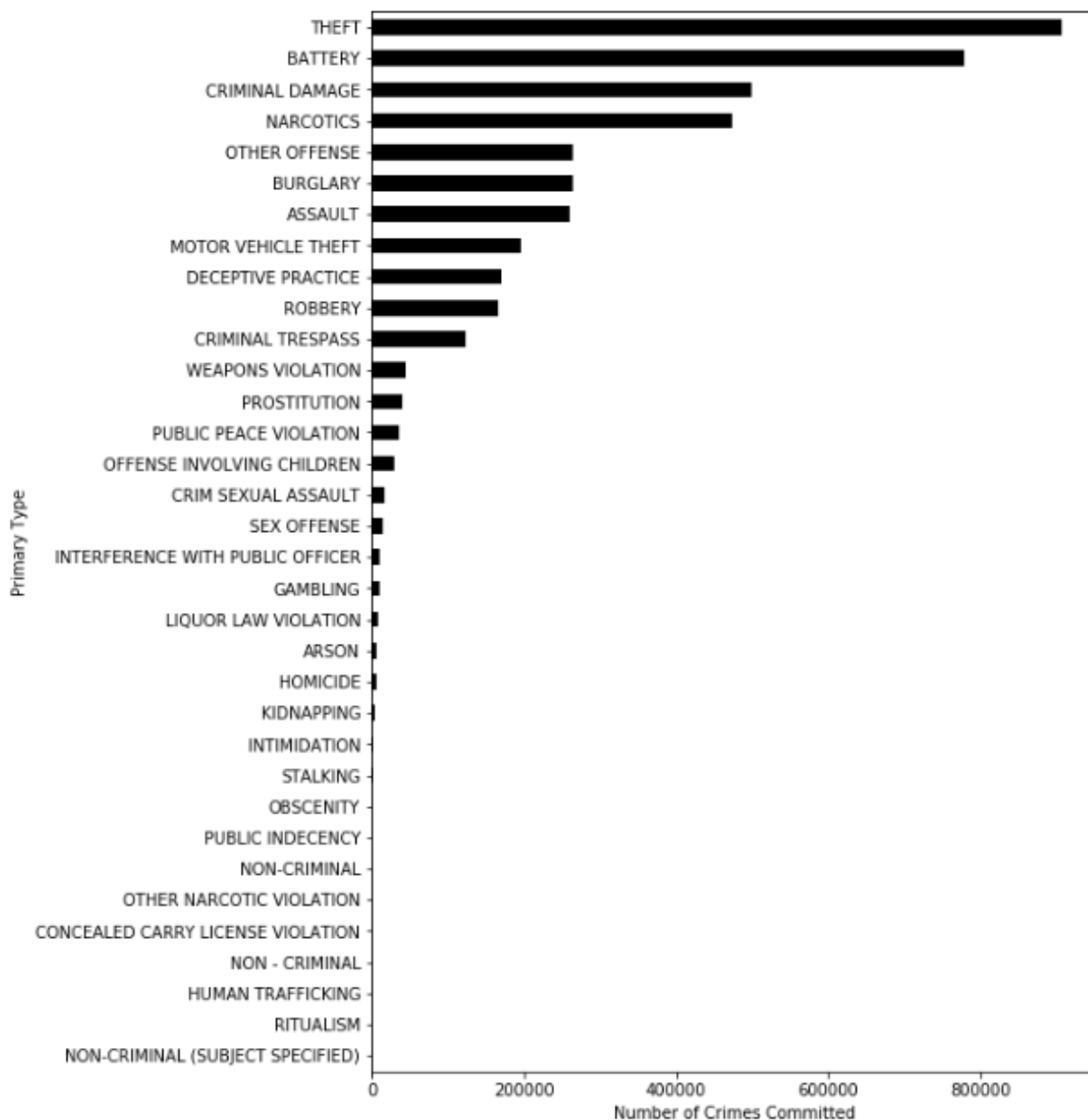


Fig. 6. Number of crimes reported for each type of crime.

The crimes for theft were the highest to be reported as high as 1,640,506, while domestic violence was the least reported crime with a number as low as 2. There appears to be a large difference between the incidents of different crimes reported. The lack of domestic violence crimes can be attributed to probable unwillingness to report domestic violence, not necessarily implying that domestic violence does not happen in the city of Chicago.



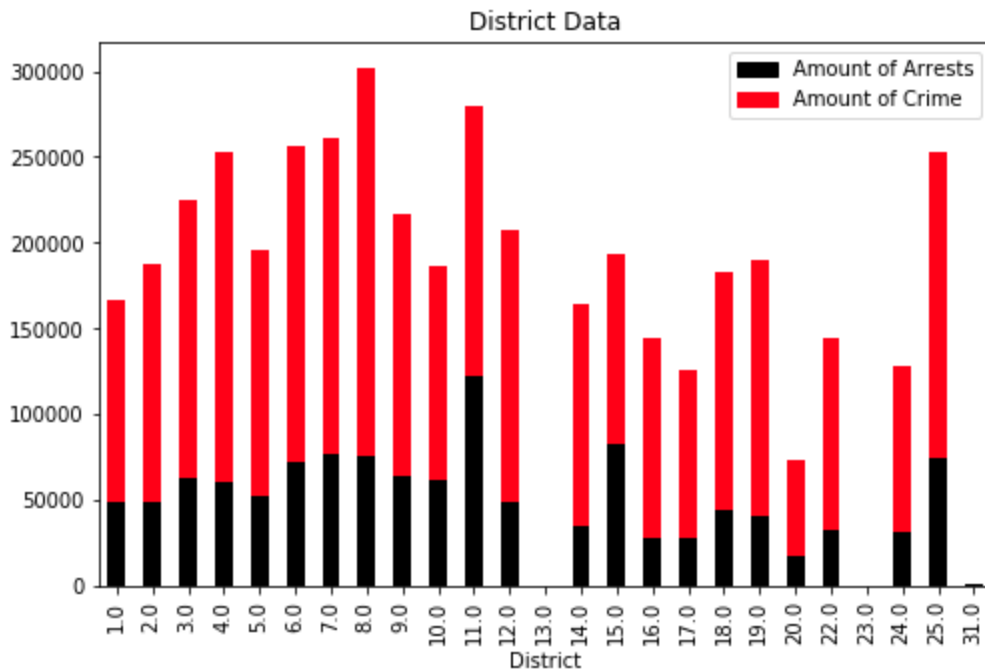


Fig. 7. Number of arrests made for the number of crimes reported district wise.

Fig. 7 shows the number of arrests that have been made in each district as compared to the number of crimes that have been committed. For District number 13, 23 and 31, there are hardly any number of crimes committed and arrests made. District 11 and 15 have a 50% arrest rate for the crimes committed.

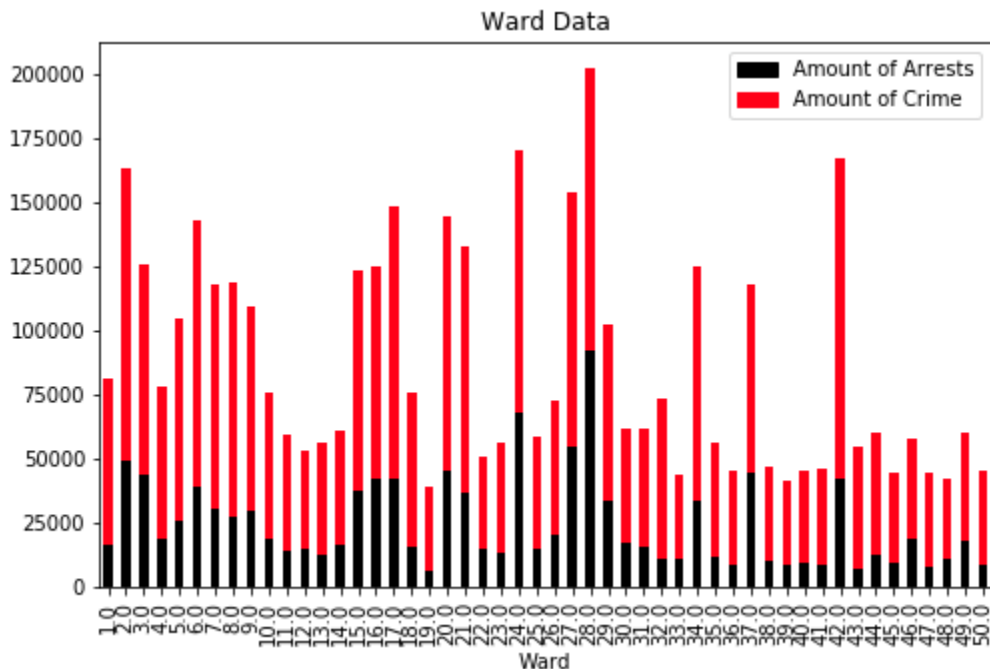


Fig. 8. Number of arrests made for the number of crimes reported ward wise.

Fig. 8 shows the number of arrests that have been made in each ward as compared to the number of crimes that have been committed. Ward numbers 28, 24 and 37 have a 50% arrest rate for the crimes committed. All other wards have about similar rates.

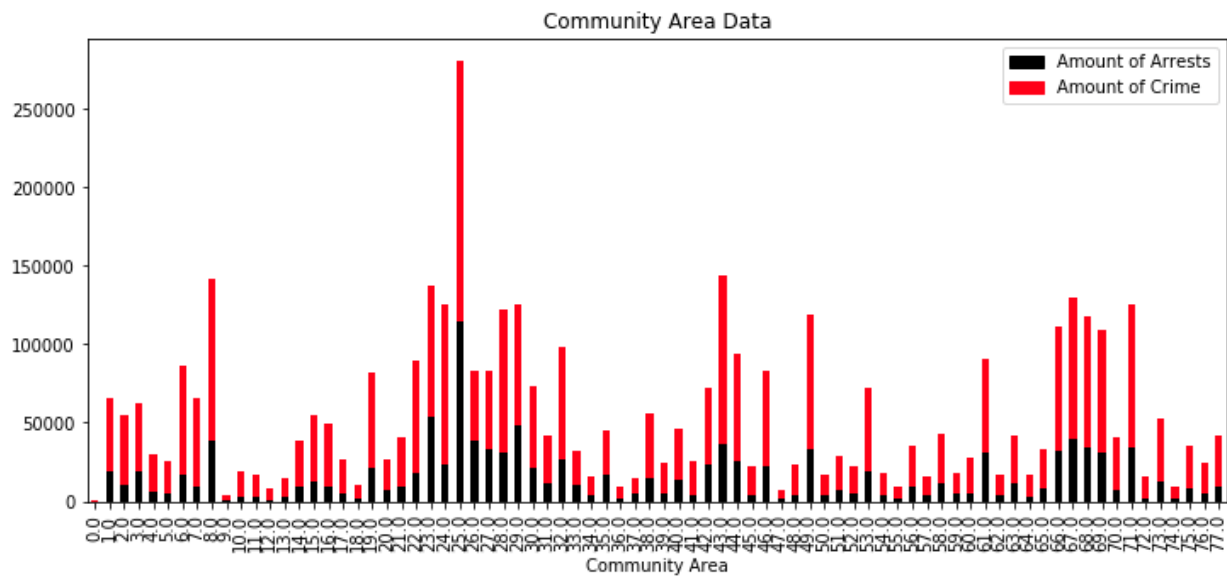


Fig. 9. Number of arrests vs the number of crimes reported community area wise.

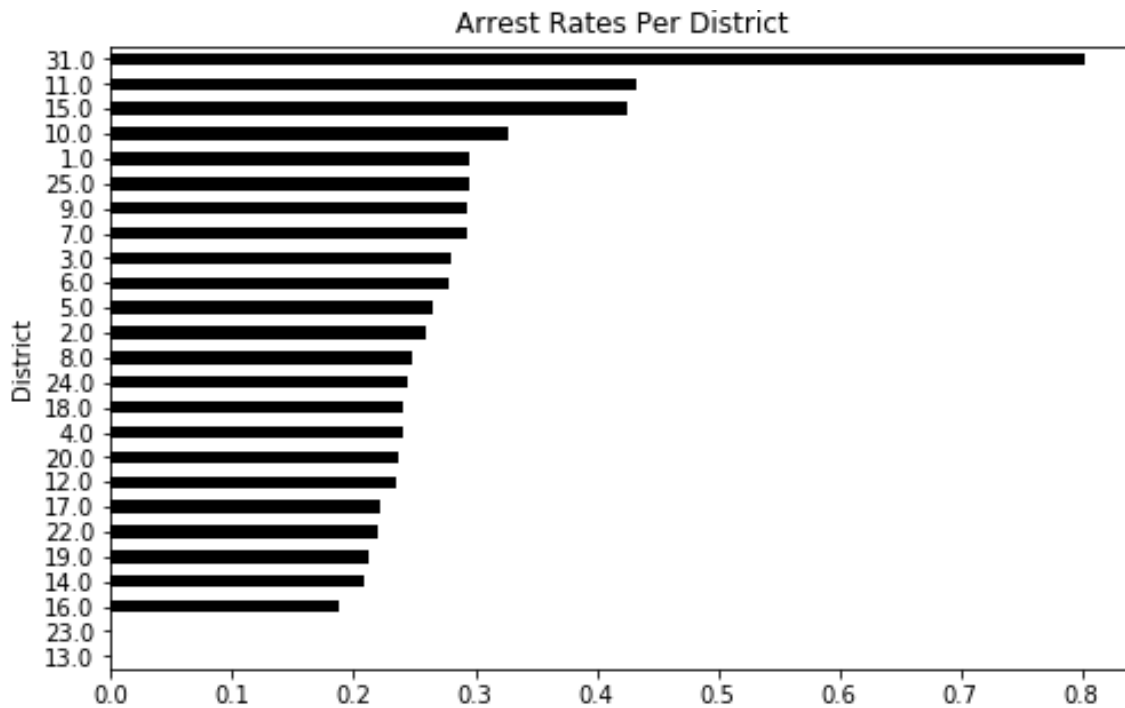


Fig. 10. Number of arrests made in districts in decreasing order.

Fig. 10 shows the arrest rates for each of the district. As seen from the data in Fig. 7, District 11 and 15 have about 50% arrest rate but District 31 has the highest arrest rate as high as 80%. This was not evident on Fig. 7 due to less number of crimes being committed in the district. District 13 and 23 have almost 0% arrest rates.

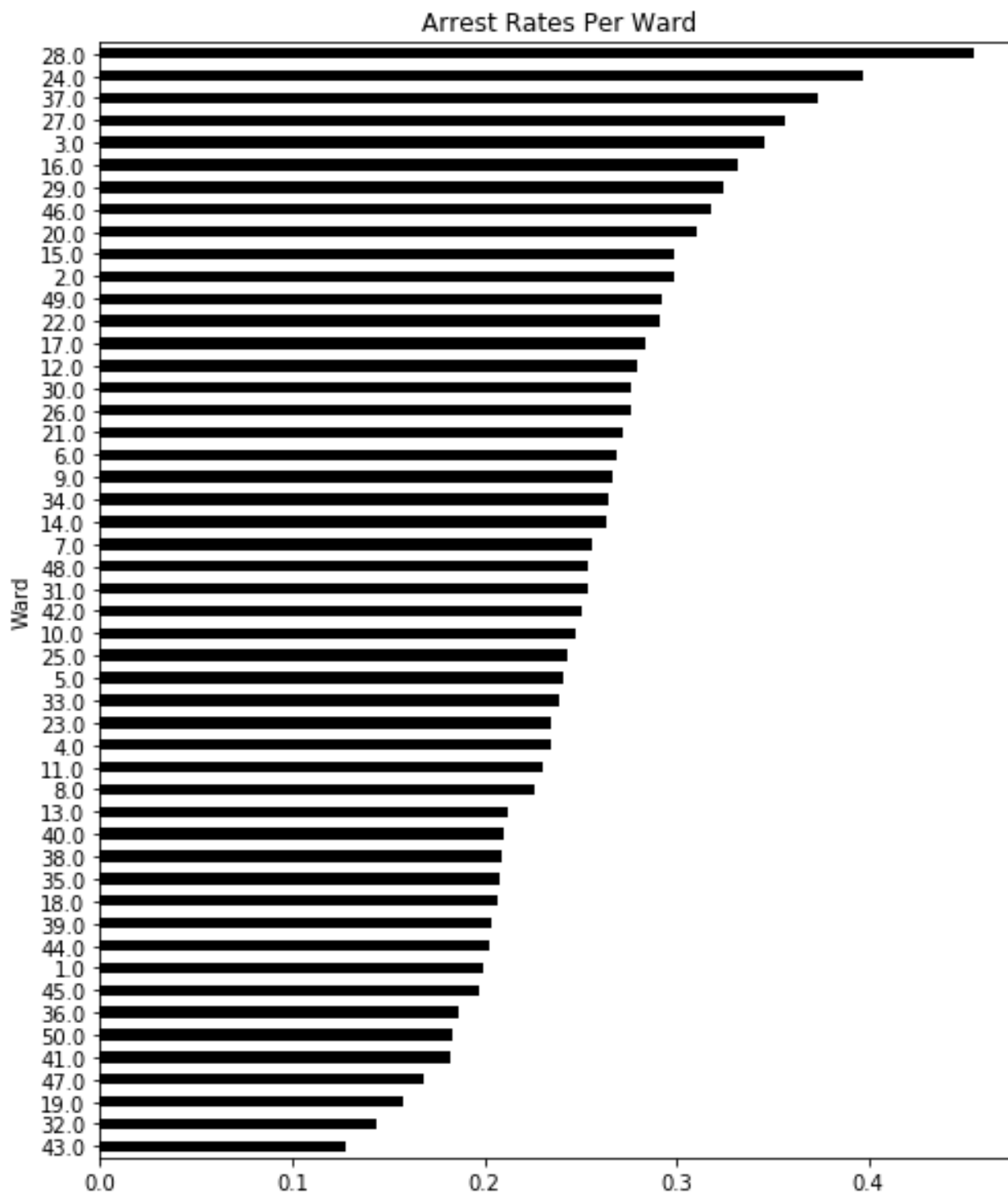


Fig. 11. Number of arrests made in wards in decreasing order.

Fig. 11 shows the arrest rates for each of the wards. As seen from the data in Fig. 8, Ward number 28 has about 50% arrest rate followed by Ward 24 and 37. On the contrary ward 43 has the lowest arrest rate as low as 12%. Rest all the wards have arrest rates between 15 and 40%.

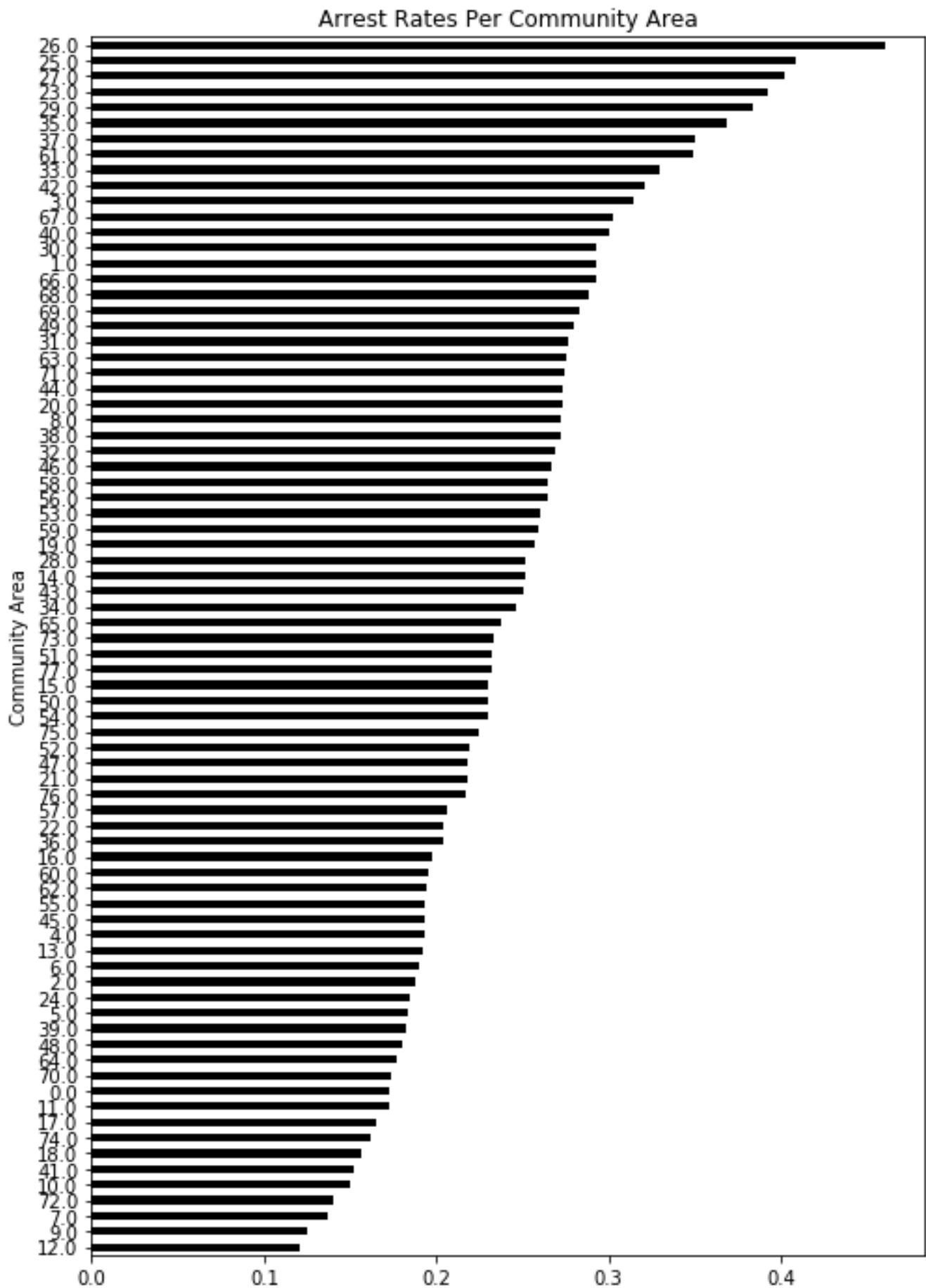


Fig. 12. Number of arrests made in wards in decreasing order.

- The number of arrests for the reported crime incidents.

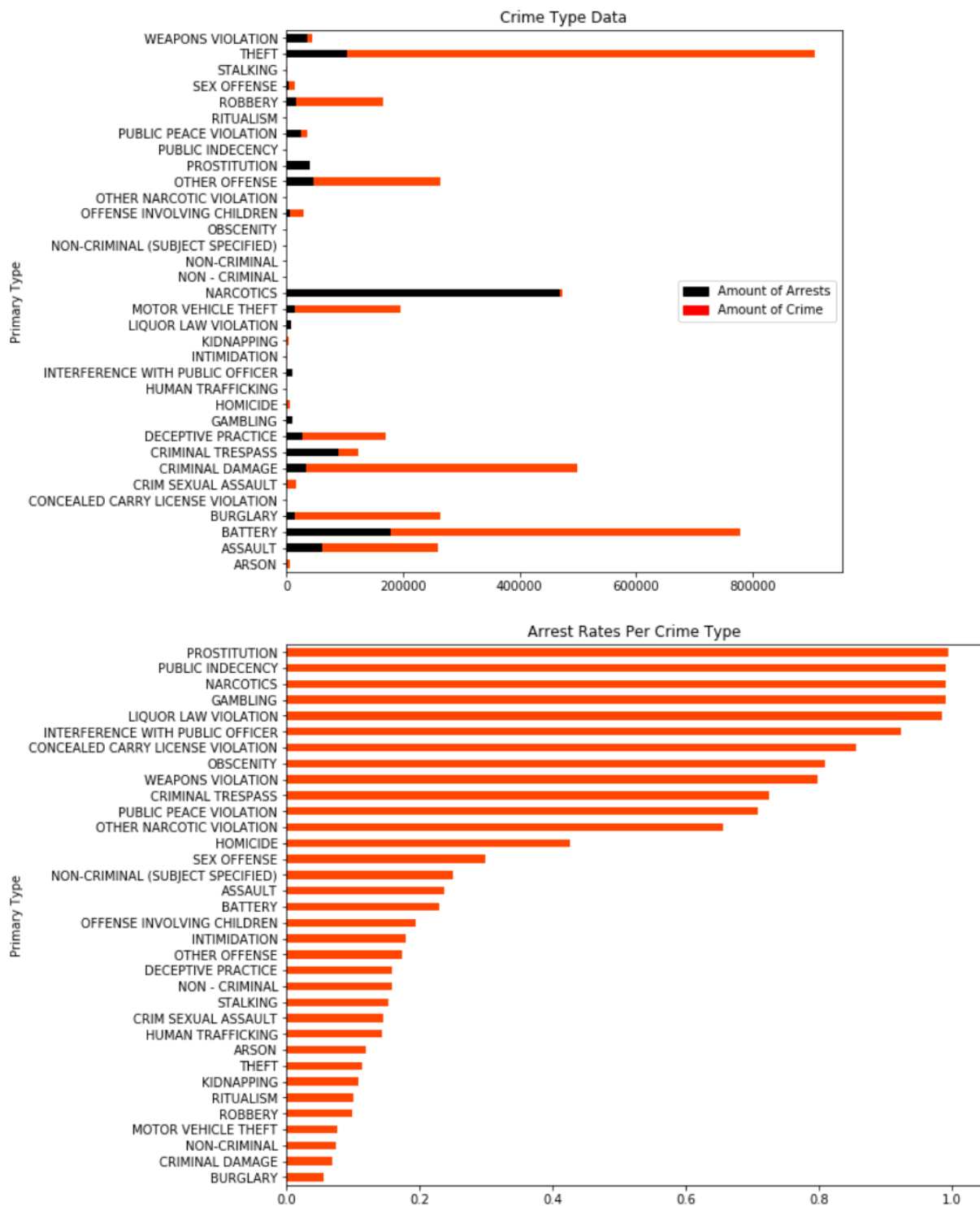


Fig. 13. Arrest vs Crime Rates

The preliminary analysis conducted provides decent insights about the crime and policing pattern in the city of Chicago. It also brings to attention certain major issues about the law enforcement and their understanding of criminal activity. We proceed to classify arrest patterns to better interpret the relationship between criminal acts and the resulting arrests.

# Explored Methods

Our general strategy when modeling the dataset was to build various classifiers, and then evaluating the results using the accuracy on the test set. Confusion matrices were generated for each model to further understand the predictive performance for the framework. The steps below are in the order that they were performed :

## 1. Principal Component Analysis

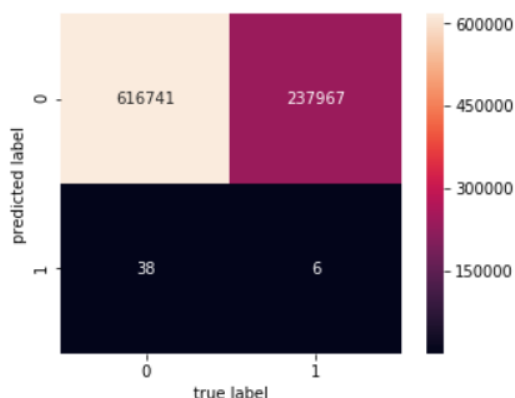
Before using the dataset on supervised learning methods, we decided to use unsupervised learning to help better explain the variance in the information. Therefore, we applied PCA to our dataset to derive a set of low-dimensional features from our variables that would potentially help us gain a clearer perspective of the relationships. We chose the number of components needed to describe the data on the basis of the cumulative *explained variance ratio* as a function of the number of components. Because we had already removed several redundant features, we observed that there was no change in the number of variables. The resulting output had 12 features which were then used in the subsequent approaches.

## 2. GLM Least Square Regression

Initially, we implemented the Linear Regression model on our dataset to predict the possibility of an arrest on the basis of the various parameters that were available. For this, we split our dataset into a training set containing 80% of our data and the remaining 20% as our test set. The probability was calculated based on the logic that if the probability of arrest is more than 50%, there would be an arrest, and for all the other cases (less than 50%), the crime would not lead to an arrest. We created a confusion matrix for our results and found that even though the test accuracy of the model was 72.15%, the model was predicting clearly wrong results that were biased. Because most of the records in our data are of the class 'not arrested', we think that the model was predicting the value on the basis of the mode of the data - classifying almost everything into the popular class. Thus, we concluded that the linear regression would not be a good model for our dataset.

	precision	recall	f1-score	support
0.0	1.00	0.72	0.84	854708
1.0	0.00	0.14	0.00	44
avg / total	1.00	0.72	0.84	854752

Accuracy of Model is: 0.721551



### 3. Logistic Regression

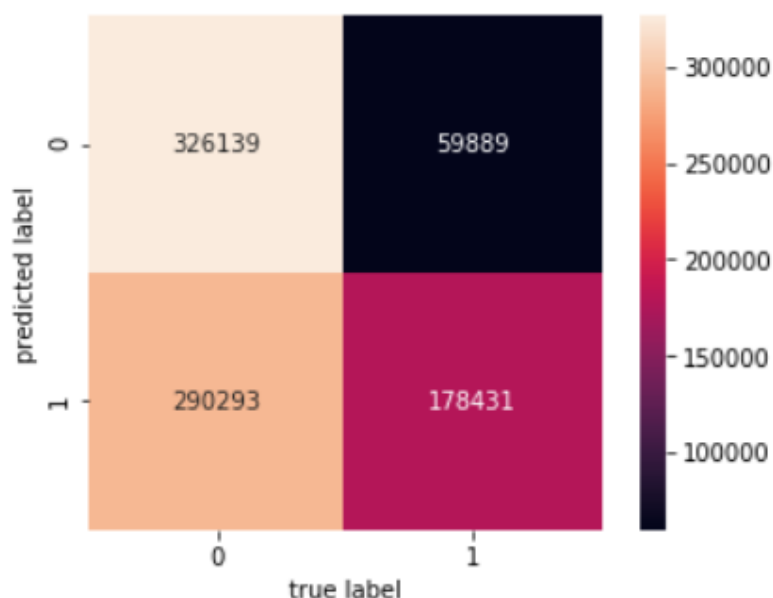
Next we considered Logistic Regression to judge whether the probability of the arrest belong to the class true or false for arrest. Again, we split our dataset into a training set containing 80% of our data and the remaining 20% as our test set. This model was only 47.88% accurate but the predictions were better than the ones we obtained from the linear regression model because it was not biased towards a false classification of arrest. As observed from the confusion matrix, this model had a relatively high number of correct classifications (false as false and true as true) and a low number of misclassifications of true arrests as false. Unfortunately, there were a high number of false arrests categorized as true which threw off the overall accuracy of the model. The team decided that although this binary classification approach was a valid framework for our data, we continued to assess other models.

```
[ 0.28300202  0.59447727  0.33211146]
Mean Cross Validation Accuracy for Logistic Regression: 0.4031969168842588
      precision    recall  f1-score   support

 False         0.53      0.84      0.65     386028
  True         0.75      0.38      0.50     468724

 avg / total         0.65      0.59      0.57     854752
```

Accuracy of Model is: 0.478898

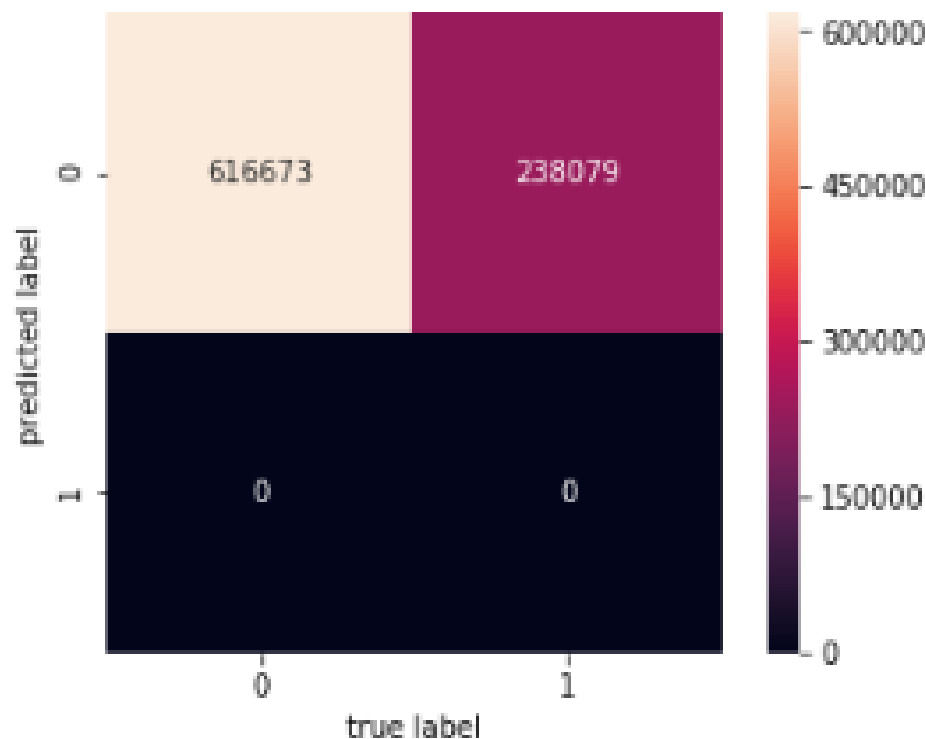


#### 4. Gaussian Naive Bayes

Next, we chose the Naive Bayes classification model by assuming features in our dataset are independent of each other - this is mostly true although it can be argued that some of the geospatial data is correlated. Once more, we split our dataset into a training set containing 80% of our data and the remaining 20% as our test set. The accuracy of the model was 72% but, as seen, the bias is even more than the Generalized Linear Model as the model classified 100% of the test set as not arrests. We believe the reason this model did not perform well was because the Gaussian classifier, assumes a normal distribution and usually works better for continuous values; on the contrary, our dataset contained mostly discrete, nominal predictors and therefore was reason to discard the model.

	precision	recall	f1-score	support
False	0.72	1.00	0.84	616673
True	0.00	0.00	0.00	238079
avg / total	0.52	0.72	0.60	854752

Accuracy of Model is: 0.721464





## 5. Bernoulli Naive Bayes

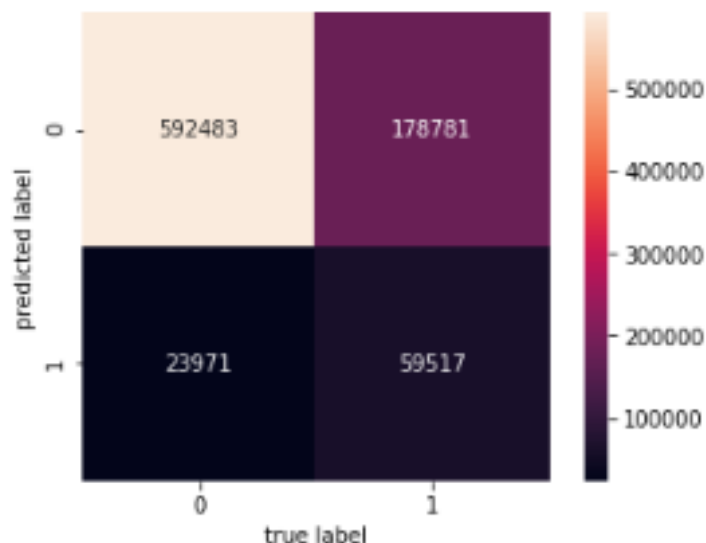
We also attempted to build Bernoulli Naive Bayes model - which works well when the variables are discrete (with the same assumption of independence). As with the previous models, we split our dataset into a training set containing 80% of our data and the remaining 20% as our test set. This model performed surprisingly well and produced an accuracy of 76.27%. Even with such a high accuracy level, when examining the confusion matrix, it can be seen that the model is still somewhat biased towards false classifications (most of the outcomes are predicted to be not arrests). That being said, this model is doing a fairly good job at not predicting false arrests as true which could be beneficial to our objective.

```
[ 0.75910913  0.7614156  0.74496713]
Mean Cross Validation Accuracy for Bernoulli Naive Bayes: 0.7551639530965487
      precision    recall  f1-score   support

 False         0.96         0.77         0.85       771264
  True         0.25         0.71         0.37         83488

 avg / total         0.89         0.76         0.81      854752

Accuracy of Model is: 0.762794
```



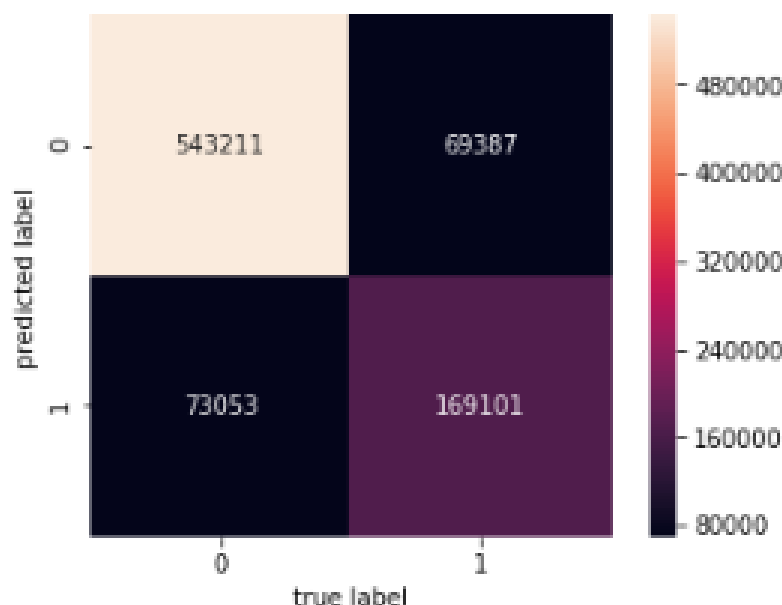
## 6. Decision Tree

We selected decision trees as our next model to predict the outcome. Again, we split our dataset into a training set containing 80% of our data and the remaining 20% as our test set. The accuracy of the model was 83.33% and produced the best results out of all the models built. Although computationally a little expensive, the model's confusion matrix was very good. The diagonal line formed by the true negatives that were classified as negative and true positives that were classified as positive were where most of the observations well. Both darker regions represent the misclassifications by the model - the fact that not a large percentage of records fell into this region suggested that this model was doing a good job at representing the relationships between the variables.

Mean Cross Validation Accuracy for Decision Tree: 0.5532466860100826

	precision	recall	f1-score	support
False	0.88	0.89	0.88	612598
True	0.71	0.70	0.70	242154
avg / total	0.83	0.83	0.83	854752

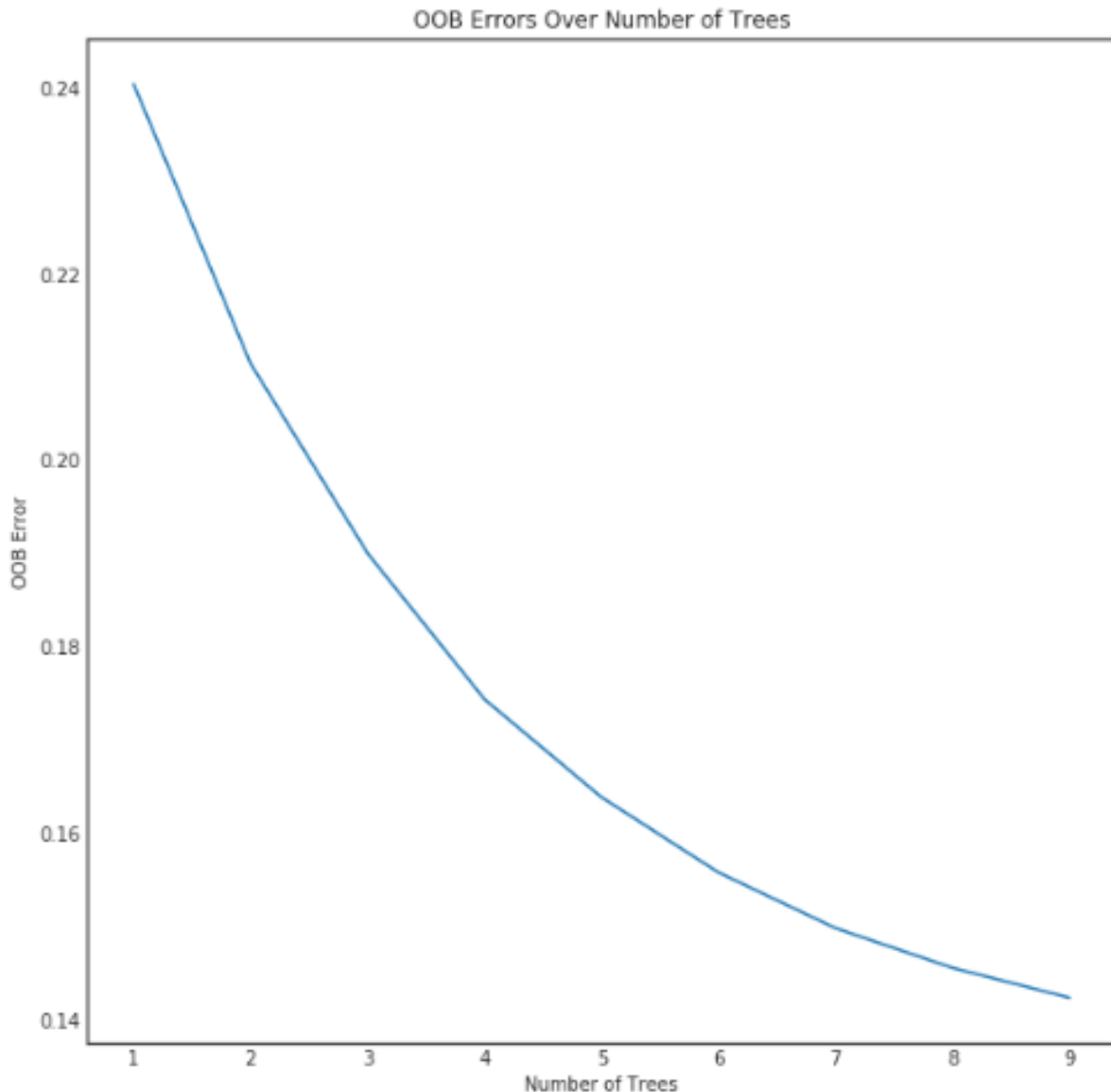
Accuracy of Model is: 0.833355



## 7. Random Forest

Given that the decision tree model performed the best out of all the previous frameworks, we next attempted to use ensemble learning to construct multiple decision trees and aggregate them into a random forest in order to produce a more optimized result. As with the previous models, we split our dataset into a training set containing 80% of our data and the remaining 20% as our test set.

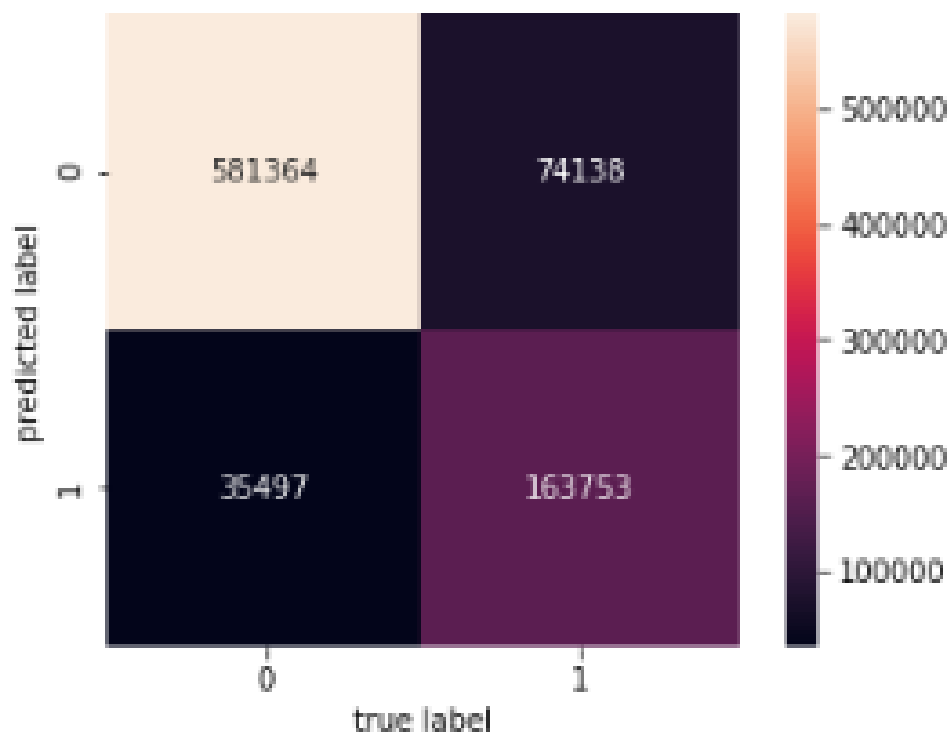
Now that we had selected a framework and were working on optimizing the model, we ran a simulation of the random forest to analyze the Out-Of-Bag (OOB) errors rates for different number of trees. As seen from the figure below, the error rate falls as more number of trees are used. This is an expected result - but a design decision was made to select 9 trees as the optimal number. Building more trees than 9 provided not much gain in accuracy as compared to the computational resources and time needed to train the model.



For the random forest with 9 trees, we got an accuracy of 87.17% - which was the highest level of correct classifications that were obtained from our experiments. As noted in the decision tree analysis, the confusion matrix displayed that most observations were being classified in the correct class.

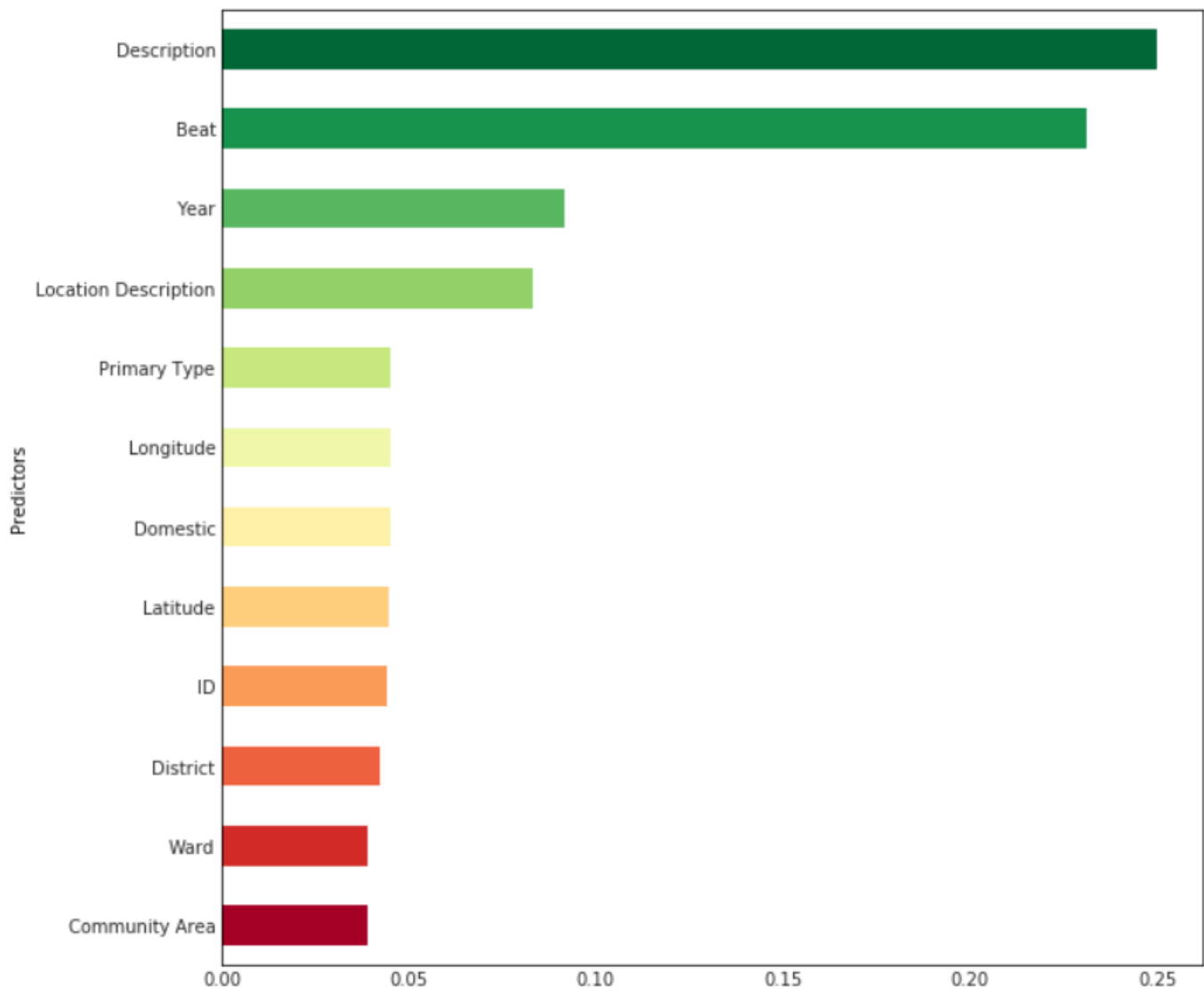
	precision	recall	f1-score	support
False	0.94	0.89	0.91	655502
True	0.69	0.82	0.75	199250
avg / total	0.88	0.87	0.88	854752

Accuracy of Model is: 0.871735



## Findings & Analysis of Models

Based on the analysis of the results we obtained from all the models, we can observe that the best model to predict the probability of arrest is the random forest. Because of its high accuracy and the correct prediction of the false arrests and the actual arrests, it was better than all the other models we implemented. As an additional benefit, the random forest and decision tree were more interpretable than some of the other models which allowed us to obtain the most important features when classifying arrest in the random forest.



We can observe that the Crime Description is the most important factor for deciding whether or not an arrest will take place. Obviously, this makes sense - the classification of arrest should be primarily based on the type of crime, as per policing protocol. We decided this was important to note as it implies that our model is considering the correct feature to calculate the possibility and produce the results.

The second most important factor was the Police Beat, which is surprising since the territory and the time of patrol of an officer should not affect the rate of arrest. After exploring the data again, we observed that the arrests were all concentrated in particular areas and were only related to certain Police Beats. This could be a potential problem for the Chicago Police Department as it could hint to many possible scenarios: performance of the other Police Beats, the possibility of some Police Beats being more aggressive than the others in their approach towards the offender, bias against certain communities, etc.

The third most important feature for the prediction was the year. As discussed in the data exploration, we concluded that as the year increased, the rate of arrest decreased as well. We can infer that the crime rate in the city of Chicago decreased in the span of 13 years and this was one of the reasons that the rate of arrest also went down. It's also interesting to note that since the year is the third most important predictor: if we hold description of crime and beat constant, when you commit the crime can decide whether or not you will be arrested. Which implies that for the same crime in the same beat, one is less likely to get arrested for a crime now than in the past as the number of arrests are decreasing.

Another implication of our analysis is that the type of location played a more important role in deciding whether there would be an arrest or not as opposed to the physical, geographical location. For example, certain location descriptions such as streets and residences had a high number of arrests but which street or where the residence is did not play as important a role.

Lastly, the features that did not show much importance on the classification of arrest were fairly expected from our analysis above. Longitude, Latitude, Community, District, Ward all had little significance on our model. The team would like to note that most of these features are correlated to other variables which might need further analysis to finalize definitive findings.

## Challenges

The biggest challenge we faced while working with our dataset was its quality. Our dataset consisted of approximately 6 million records and had inconsistencies and missing data. This limited our scope of working on questions we had planned on working on, in the initial stages of the project. We did initial work on answering two other questions to give better insights to Chicago Police Department which are as below:

- What better security measures can be implemented in which areas?

We wanted to do a detailed analysis to identify if there were any specific type of crime that had high frequencies in some areas. On the basis of the location where these crimes were occurring, suggestions could be provided as to what type of security measures could be implemented in that district or block to avoid the high crime rate. For instance, if there are lots of mugging reported on a particular block, installing more surveillance cameras can help in reducing the number. Since the reported time was inconsistent and the addresses were masked, we did not reach to substantial conclusions.

- How to better assess incomplete 911 calls?

We wanted to predict the type of crime that should be expected from a particular block or ward if the 911 caller is unable to give all the information. That would help the EMS team to estimate the amount of backup required.

Further, most of our data was either geographical data, date time and categorical data, which limited the implementation of models to only qualitative models. Because of the large size of the dataset and limited computational power, running the models was time consuming task. Some of the models took up to 10 hours of running time limiting our scope of exploration.

## Weaknesses

The challenges discussed previously led to most of the weaknesses in our modeling and analysis. These areas could have led to inaccurate results and the team believed it was important to note these down:

- The geographical nature of the data, with multiple factors playing a role in the frequency and occurrence of crime, abstained the use of a general model for the entire data while predicting results.
- Lack of certain demographics (like unemployment, age, gender, alcoholism etc.) left dependent attributes from the database to be inadequate, which may have led to biased predictions (collinearity).

- Addresses have been masked to protect the privacy of crime victims, aggregating addresses to block level. This affected the accuracy of predictions.
- Computational and Time Resource available to us were limited. Therefore, a balance needed to be met between model complexity and accuracy.

## Conclusion

An analysis on the relationships between crime and arrests for the city of Chicago can help city authorities plan, prepare, organize, and allocate available security measures to improve safety concerns. Our models revealed that, for the most part, the police department is doing a commendable job in assessing the classification of an arrest based solely on the crime itself. On the other hand, several predictors that were surprisingly impactful on the arrests proved to be the beat and the year of the crime. The team recommends a further, in depth analysis on the arrest rate related to these variables. A more thorough understanding of the systematic reduction in arrests over time would provide a clearer snapshot of the criminal activity in the city. It's also imperative some investigation behind the distribution of arrest rates among the various beats is conducted to ensure the policing arounds these crimes is fair process. Before implementing any changes in policies that address these concerns we hope to continue analyzing the information and overcome the challenges and weaknesses we faced in our analysis to assist the city in their pursuit towards a safe future.



## References

1. (What's Behind Chicago's Surge in Violence). Retrieved March 20, 2018, from <http://www.chicagotribune.com/news/opinion/editorials/ct-chicago-crime-increase-causes-edit-0118-md-20170117-story.html>
2. (Chicago Crime Data). Retrieved March 02, 2018, from (Chi<https://www.kaggle.com/currie32/crimes-in-chicago/data>
3. (An Introduction to Statistical Learning). Retrieved March 02, 2018, from <http://www-bcf.usc.edu/~gareth/ISL/>
4. Boundaries - Police Beats (current) | City of Chicago | Data Portal. (n.d.). Retrieved March 02, 2018, from <https://data.cityofchicago.org/d/aerh-rz74>
5. C. (2017, January 28). Retrieved March 02, 2018, from <https://www.kaggle.com/currie32/crimes-in-chicago/data>
6. Chicago Police Department - Illinois Uniform Crime Reporting (IUCR) Codes | City of Chicago | Data Portal. (n.d.). Retrieved March 02, 2018, from <https://data.cityofchicago.org/Public-Safety/Chicago-Police-Department-Illinois-Uniform-Crime-R/c7ck-438e>
7. (Crime Type Categories). Retrieved March 02, 2018, from [http://gis.chicagopolice.org/clearmap\\_crime\\_sums/crime\\_types.html](http://gis.chicagopolice.org/clearmap_crime_sums/crime_types.html)
8. (Google Maps APIs). Retrieved March 02, 2018, from <https://developers.google.com/maps/>
9. (Chicago Crime Analysis). Retrieved March 02, 2018, from <https://github.com/ajitkoduri/Chicago-Crime-Analysis>
10. (Public Safety). Retrieved March 02, 2018, from <https://github.com/dssg/publicsafety>

## Appendix

Find Python code used to work with the data below