# FOUNDATION OF DATA SCIENCE
# CSGY-6053

# PROJECT REPORT:
# FORECASTING NYC CRIME DATA

Harsh Sonthalia – hs4226
Jacob Fishman – jf4322
Ujjwal Vikram Kulkarni – uk2011

# Introduction

We have all moved from far and wide to live in the greatest city in the world. With all of the excitement from moving comes worries about safety, more importantly in terms of crime. There are plenty of statistics that show that crime in New York city is concerning, from the elevated violent crime rate to crimes per square mile in New York City totalling 1,549 compared to the national average of 28.3.[1] The goal of this project is to model the time series of crime trends in New York City, and forecast the crimes for the next year, classifying these predictions to better help with preventive measures for the police department and local citizens.

# Pre Processing

We are using the data set from NYC Open Data Website:

https://data.cityofnewyork.us/Public-Safety/NYPD-Complaint-Data-Historic/qgea-i56i

This dataset provides 7.38 Million instances of crimes, all of the crimes that have been reported from 2006 onwards, with 35 features. We have gone through the dataset and removed any duplicate or unnecessary features, which has reduced our features to 28. For the rest of the preprocessing we will use the data exploration to determine the factors that are relevant to our predictions.
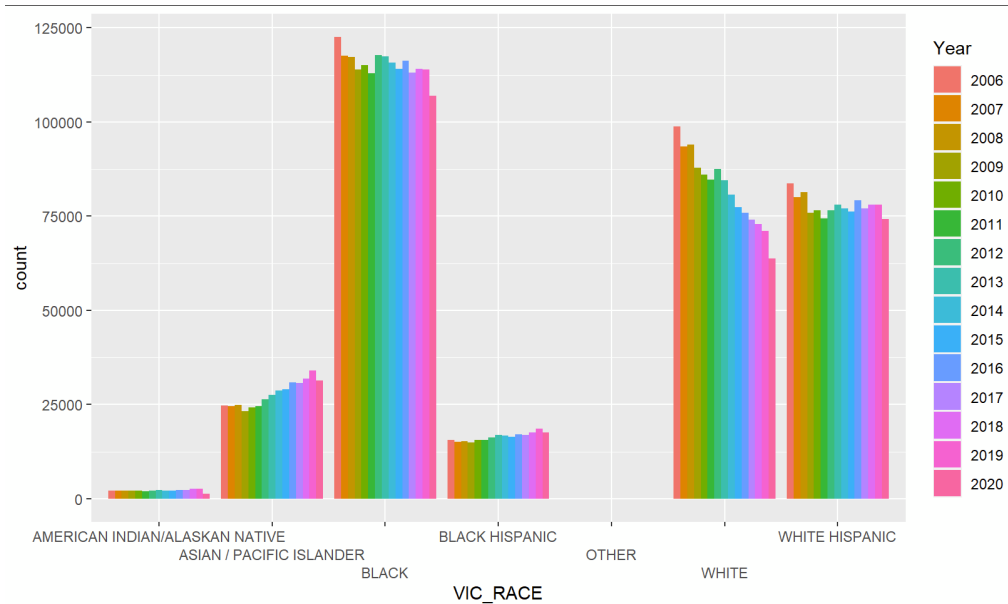
# Data Exploration

1) Here is the Pearson correlation matrix for all of our features that we deemed to be helpful in understanding crime data. You can see that there are no features that are strongly correlated. The highest absolute value is between Key ID and Law Category ID, both of which should be correlated since they deal with the classification of the crime itself.
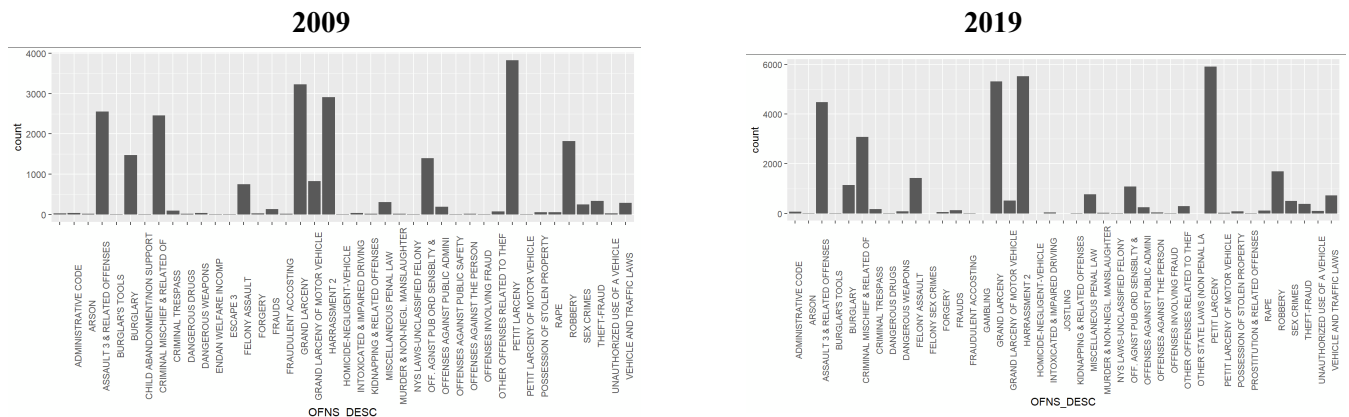
Out[37]:

|  | KY_CD | LAW_CAT_CD | BORO_NM | LOC_OF_OCCUR_DESC | PREM_TYP_DESC | JURISDICTION_CODE | VIC_AGE_GROUP | VIC_RACE | VIC_SEX | Start_Year | Start_Month | Start_Day | Start_Hour | Latitude | Longitude |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| KY_CD | 1.000000 | 0.562541 | 0.021159 | -0.006615 | -0.026598 | -0.000648 | -0.008442 | -0.061142 | -0.110393 | 0.029934 | -0.011171 | -0.000409 | 0.055742 | 0.001667 | -0.020599 |
| LAW_CAT_CD | 0.562541 | 1.000000 | 0.006622 | -0.020458 | -0.028169 | 0.004113 | -0.029517 | -0.042883 | -0.035857 | -0.007274 | -0.011854 | 0.000756 | 0.033861 | 0.014443 | -0.015889 |
| BORO_NM | 0.021159 | 0.006622 | 1.000000 | -0.057136 | -0.079594 | -0.021906 | 0.023183 | 0.046668 | 0.001188 | -0.013862 | -0.001459 | -0.002351 | 0.011035 | -0.732659 | -0.156119 |
| LOC_OF_OCCUR_DESC | -0.006615 | -0.020458 | -0.057136 | 1.000000 | 0.112384 | 0.017863 | 0.009262 | -0.028285 | -0.112284 | -0.009906 | -0.020942 | -0.007321 | -0.023090 | 0.047809 | 0.003901 |
| PREM_TYP_DESC | -0.026598 | -0.028169 | -0.079594 | 0.112384 | 1.000000 | 0.018256 | 0.015188 | 0.059688 | 0.014440 | 0.010174 | -0.007152 | -0.001578 | -0.026262 | -0.003255 | -0.063768 |
| JURISDICTION_CODE | -0.000648 | 0.004113 | -0.021906 | 0.017863 | 0.018256 | 1.000000 | 0.001235 | -0.022001 | 0.001737 | 0.009150 | -0.002017 | 0.000144 | -0.000967 | 0.011098 | -0.007135 |
| VIC_AGE_GROUP | -0.008442 | -0.029517 | 0.023183 | 0.009262 | 0.015188 | 0.001235 | 1.000000 | 0.021549 | 0.050705 | 0.002629 | -0.002380 | -0.010874 | 0.025085 | -0.023415 | 0.008193 |
| VIC_RACE | -0.061142 | -0.042883 | 0.046668 | -0.028285 | 0.059688 | -0.022001 | 0.021549 | 1.000000 | 0.101038 | -0.004089 | 0.009573 | -0.000360 | 0.011485 | -0.186903 | -0.129119 |
| VIC_SEX | -0.110393 | -0.035857 | 0.001188 | -0.112284 | 0.014440 | 0.001737 | 0.050705 | 0.101038 | 1.000000 | 0.015806 | 0.010040 | 0.004193 | -0.010907 | -0.012719 | -0.012062 |
| Start_Year | 0.029934 | -0.007274 | -0.013862 | -0.009906 | 0.010174 | 0.009150 | 0.002629 | -0.004089 | 0.015806 | 1.000000 | 0.001693 | 0.001008 | -0.001196 | 0.026819 | 0.018626 |
| Start_Month | -0.011171 | -0.011854 | -0.001459 | -0.020942 | -0.007152 | -0.002017 | -0.002380 | 0.009573 | 0.010040 | 0.001693 | 1.000000 | 0.007017 | 0.002566 | -0.000272 | -0.001152 |
| Start_Day | -0.000409 | 0.000756 | -0.002351 | -0.007321 | -0.001578 | 0.000144 | -0.010874 | -0.000360 | 0.004193 | 0.001008 | 0.007017 | 1.000000 | 0.014746 | 0.002929 | 0.001548 |
| Start_Hour | 0.055742 | 0.033861 | 0.011035 | -0.023090 | -0.026262 | -0.000967 | 0.025085 | 0.011485 | -0.010907 | -0.001196 | 0.002566 | 0.014746 | 1.000000 | -0.005675 | 0.005380 |
| Latitude | 0.001667 | 0.014443 | -0.732659 | 0.047809 | -0.003255 | 0.011098 | -0.023415 | -0.186903 | -0.012719 | 0.026819 | -0.000272 | 0.002929 | -0.005675 | 1.000000 | 0.300621 |
| Longitude | -0.020599 | -0.015889 | -0.156119 | 0.003901 | -0.063768 | -0.007135 | 0.008193 | -0.129119 | -0.012062 | 0.018626 | -0.001152 | 0.001548 | 0.005380 | 0.300621 | 1.000000 |

2) The next figure shows the relation between total crimes and the victim's reported race, grouped by year. Here we can see that there is a declining crime rate for most of the race groups, which follows the trend we see with the crimes per year. However the "Asian/Pacific Islander" group have seen an increase in being victims in crimes:



We evaluated the data for the years 2009 (the lowest rate) and 2019 (the highest Rate) to dig in deeper regarding this upward trend:



The largest increase in total crimes in specific offense categories from 2009 to 2019 is Harassment 3 and Assault 3, both of which are physical crimes. These results indicate that victim race would be an important feature in determining classification of crimes.

By looking at the total crimes per year, we see a downward trend from 2007 to 2020. This trend will be useful to address in terms of our forecasting.



Yearly distribution of crime in NYC

Even though the number of crimes have been decreasing since 2008, the rate of change has not been constant and does not follow any evident pattern. This indicates we do not need to use the year as a feature in our KNN Model. The sudden drop in 2020 is obvious because of the restrictions imposed due to COVID-19.



TOTAL CRIME COUNT DIFFERENCE

Prior to our analysis we expected the preponderance of crimes to occur late at night. Surprisingly the data shows that crime rates per hour gradually increases from 5AM and reaches its peak in the evenings.

**Hourly distribution of crime in NYC**



This is an intriguing plot as it provides a clear indication that the number of crimes on the 1st day of the month is the highest when compared to the rest of the days.

**Heatmaps:**

Year vs Month



- There is a higher frequency of crimes from May to October most likely because more people venture out in the summer than in winters.
- April 2020 had the lowest crime count seemingly due to the COVID Lockdown.

## Month vs Day of the week



- Most crimes happen on Friday compared to Sunday where the crime count is the lowest.

## Hours vs Day of the week

| hour | Friday | Monday | Saturday | Sunday | Thursday | Tuesday | Wednesday |
|------|--------|--------|----------|--------|----------|---------|-----------|
| 0.0 | 47430.0 | 41744.0 | 55984.0 | 52869.0 | 43464.0 | 39446.0 | 42975.0 |
| 1.0 | 31531.0 | 26195.0 | 47169.0 | 47031.0 | 28113.0 | 24260.0 | 27423.0 |
| 2.0 | 24851.0 | 20104.0 | 41094.0 | 41333.0 | 21718.0 | 18556.0 | 20633.0 |
| 3.0 | 20149.0 | 16270.0 | 36370.0 | 38512.0 | 16864.0 | 14408.0 | 16044.0 |
| 4.0 | 16811.0 | 14784.0 | 32546.0 | 35257.0 | 13952.0 | 12411.0 | 13155.0 |
| 5.0 | 12422.0 | 11934.0 | 20076.0 | 21863.0 | 11268.0 | 10575.0 | 10821.0 |
| 6.0 | 15289.0 | 13754.0 | 15049.0 | 15006.0 | 14691.0 | 13763.0 | 14608.0 |
| 7.0 | 22491.0 | 22182.0 | 15513.0 | 14592.0 | 23016.0 | 22529.0 | 23280.0 |
| 8.0 | 37000.0 | 37195.0 | 23376.0 | 20878.0 | 36364.0 | 37123.0 | 37869.0 |
| 9.0 | 39476.0 | 40078.0 | 28824.0 | 24574.0 | 39570.0 | 40195.0 | 40127.0 |
| 10.0 | 40289.0 | 41624.0 | 33577.0 | 28467.0 | 40715.0 | 41150.0 | 40751.0 |
| 11.0 | 41213.0 | 39487.0 | 35265.0 | 30473.0 | 42130.0 | 42044.0 | 42750.0 |
| 12.0 | 58086.0 | 55437.0 | 49226.0 | 42141.0 | 57452.0 | 57712.0 | 58927.0 |
| 13.0 | 49614.0 | 45150.0 | 42046.0 | 36110.0 | 50087.0 | 50088.0 | 51551.0 |
| 14.0 | 57568.0 | 50797.0 | 45177.0 | 39655.0 | 56000.0 | 56497.0 | 59014.0 |
| 15.0 | 65405.0 | 58842.0 | 49569.0 | 44263.0 | 63088.0 | 63946.0 | 65407.0 |
| 16.0 | 62654.0 | 54438.0 | 50495.0 | 45181.0 | 59721.0 | 60937.0 | 61824.0 |
| 17.0 | 63668.0 | 54948.0 | 51647.0 | 46750.0 | 59184.0 | 60137.0 | 61310.0 |
| 18.0 | 63691.0 | 55458.0 | 53223.0 | 48422.0 | 59748.0 | 61010.0 | 62242.0 |
| 19.0 | 60364.0 | 52083.0 | 52430.0 | 46898.0 | 57463.0 | 58308.0 | 59324.0 |
| 20.0 | 59847.0 | 50780.0 | 54385.0 | 47887.0 | 56221.0 | 56495.0 | 56821.0 |
| 21.0 | 55395.0 | 44188.0 | 51362.0 | 43856.0 | 50128.0 | 49931.0 | 50916.0 |
| 22.0 | 54564.0 | 40446.0 | 51140.0 | 41739.0 | 46701.0 | 45866.0 | 45919.0 |
| 23.0 | 53280.0 | 35393.0 | 50804.0 | 37419.0 | 42328.0 | 40350.0 | 40622.0 |

day_name

- This chart, combining daily and hourly crime counts, indicates more crimes occurring when we would expect more people to be out at night, which is dusk to night on the weekends.

# Days vs Months

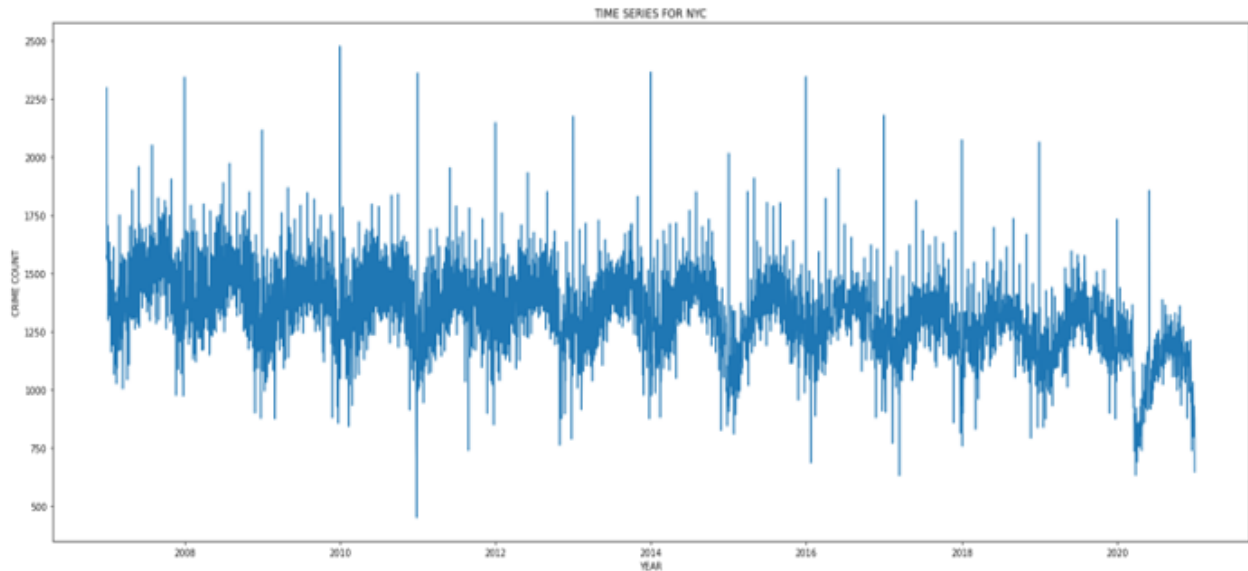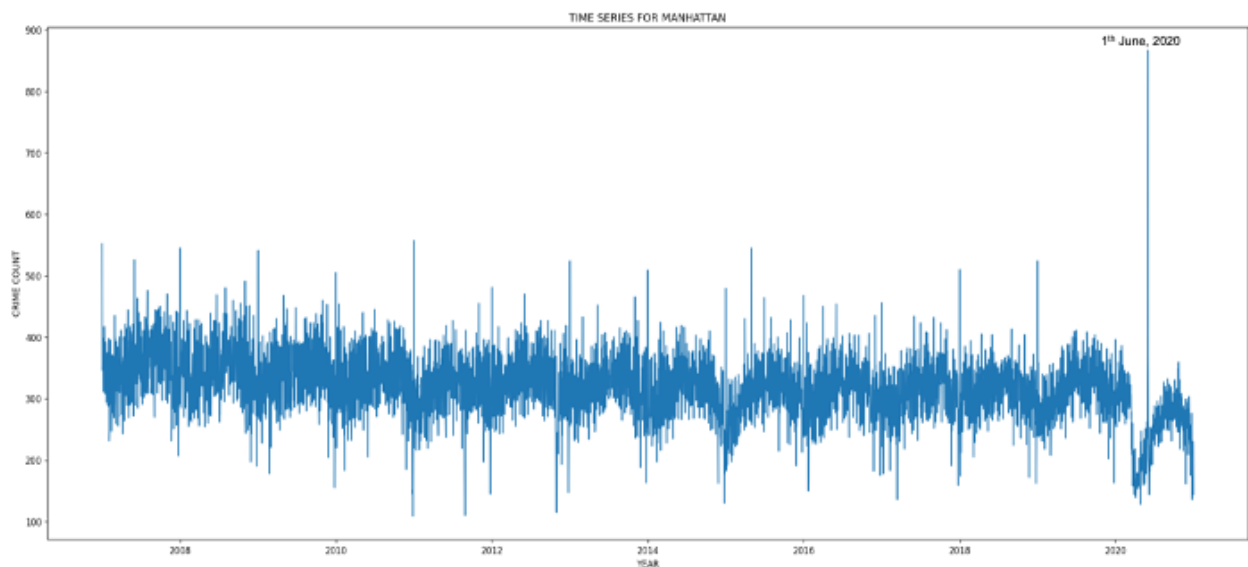| day | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 30692.0 | 21530.0 | 21907.0 | 21834.0 | 22827.0 | 24947.0 | 23493.0 | 23885.0 | 23476.0 | 22866.0 | 22418.0 | 21254.0 |
| 2 | 16760.0 | 18181.0 | 17372.0 | 18192.0 | 19327.0 | 19792.0 | 19242.0 | 19900.0 | 19418.0 | 19335.0 | 18291.0 | 17930.0 |
| 3 | 17087.0 | 17167.0 | 17605.0 | 17912.0 | 18814.0 | 19038.0 | 19918.0 | 19712.0 | 19135.0 | 19371.0 | 18004.0 | 18223.0 |
| 4 | 17000.0 | 17498.0 | 17457.0 | 17779.0 | 18535.0 | 18991.0 | 19190.0 | 19333.0 | 19756.0 | 19373.0 | 18074.0 | 18138.0 |
| 5 | 17936.0 | 17439.0 | 18167.0 | 17598.0 | 18846.0 | 19342.0 | 20287.0 | 19881.0 | 20082.0 | 19737.0 | 18300.0 | 18314.0 |
| 6 | 17285.0 | 17582.0 | 17642.0 | 17430.0 | 19133.0 | 19280.0 | 18894.0 | 19531.0 | 18608.0 | 19594.0 | 17925.0 | 17513.0 |
| 7 | 17002.0 | 17662.0 | 17820.0 | 17752.0 | 19025.0 | 18662.0 | 19288.0 | 19732.0 | 19071.0 | 19099.0 | 18586.0 | 17750.0 |
| 8 | 17609.0 | 17390.0 | 18170.0 | 17608.0 | 18861.0 | 19112.0 | 18925.0 | 19900.0 | 19561.0 | 18903.0 | 17926.0 | 17510.0 |
| 9 | 17838.0 | 16577.0 | 18484.0 | 17859.0 | 18865.0 | 18611.0 | 19275.0 | 19658.0 | 19048.0 | 19110.0 | 18437.0 | 17192.0 |
| 10 | 18495.0 | 17240.0 | 18650.0 | 18403.0 | 19673.0 | 19301.0 | 19863.0 | 20132.0 | 19324.0 | 19993.0 | 18609.0 | 18385.0 |
| 11 | 17986.0 | 16992.0 | 18255.0 | 18016.0 | 19217.0 | 19206.0 | 19204.0 | 19649.0 | 19315.0 | 19080.0 | 17489.0 | 17907.0 |
| 12 | 18153.0 | 17500.0 | 18649.0 | 17910.0 | 19071.0 | 19916.0 | 19712.0 | 19743.0 | 19697.0 | 18876.0 | 18001.0 | 18310.0 |
| 13 | 18228.0 | 16858.0 | 18123.0 | 18457.0 | 18355.0 | 19227.0 | 19362.0 | 19728.0 | 19096.0 | 18890.0 | 18196.0 | 17686.0 |
| 14 | 18116.0 | 16922.0 | 17774.0 | 18089.0 | 19076.0 | 19365.0 | 19512.0 | 19990.0 | 19527.0 | 18881.0 | 18307.0 | 17599.0 |
| 15 | 19010.0 | 18114.0 | 19083.0 | 18524.0 | 20079.0 | 20406.0 | 20180.0 | 20667.0 | 20561.0 | 19850.0 | 18682.0 | 18194.0 |
| 16 | 18090.0 | 17042.0 | 17494.0 | 18000.0 | 19219.0 | 19557.0 | 19750.0 | 19724.0 | 19225.0 | 18903.0 | 18130.0 | 16875.0 |
| 17 | 18188.0 | 17134.0 | 17607.0 | 18510.0 | 19384.0 | 19306.0 | 19802.0 | 19970.0 | 19585.0 | 19573.0 | 18145.0 | 17015.0 |
| 18 | 17736.0 | 17200.0 | 18231.0 | 18977.0 | 19003.0 | 19133.0 | 19487.0 | 18983.0 | 19757.0 | 19108.0 | 18424.0 | 17511.0 |
| 19 | 17914.0 | 16736.0 | 17570.0 | 18073.0 | 18258.0 | 19351.0 | 19591.0 | 19173.0 | 19478.0 | 18756.0 | 18227.0 | 16940.0 |
| 20 | 18285.0 | 17681.0 | 18747.0 | 18635.0 | 19978.0 | 20273.0 | 20071.0 | 20396.0 | 20135.0 | 20053.0 | 19202.0 | 17851.0 |
| 21 | 16682.0 | 17433.0 | 17580.0 | 18093.0 | 19609.0 | 19498.0 | 19533.0 | 19378.0 | 19696.0 | 19513.0 | 18496.0 | 17458.0 |
| 22 | 17547.0 | 17218.0 | 17864.0 | 18135.0 | 19391.0 | 19705.0 | 19558.0 | 19494.0 | 19487.0 | 19004.0 | 16903.0 | 17403.0 |
| 23 | 17536.0 | 16981.0 | 17629.0 | 18384.0 | 19415.0 | 19419.0 | 19429.0 | 19495.0 | 19435.0 | 19633.0 | 17350.0 | 17384.0 |
| 24 | 17139.0 | 17049.0 | 17706.0 | 18361.0 | 18927.0 | 19564.0 | 20003.0 | 19846.0 | 19356.0 | 19119.0 | 16736.0 | 15175.0 |
| 25 | 17429.0 | 17035.0 | 17956.0 | 18353.0 | 19309.0 | 19922.0 | 20042.0 | 19459.0 | 19316.0 | 19033.0 | 17568.0 | 12056.0 |
| 26 | 17224.0 | 16973.0 | 17573.0 | 17882.0 | 18810.0 | 20018.0 | 19668.0 | 19283.0 | 19280.0 | 18763.0 | 16706.0 | 14820.0 |
| 27 | 17020.0 | 18090.0 | 17689.0 | 18393.0 | 19092.0 | 19558.0 | 19730.0 | 18823.0 | 19517.0 | 18209.0 | 17143.0 | 15530.0 |
| 28 | 17091.0 | 17897.0 | 18210.0 | 18194.0 | 19035.0 | 19331.0 | 19890.0 | 19013.0 | 19524.0 | 18282.0 | 17524.0 | 16177.0 |
| 29 | 17058.0 | 4956.0 | 16934.0 | 18278.0 | 19238.0 | 19122.0 | 19574.0 | 19125.0 | 19223.0 | 17689.0 | 17341.0 | 15743.0 |
| 30 | 17294.0 | 0.0 | 17416.0 | 18375.0 | 19390.0 | 19175.0 | 19345.0 | 19656.0 | 18780.0 | 18347.0 | 18118.0 | 15923.0 |
| 31 | 17069.0 | 0.0 | 16916.0 | 0.0 | 18870.0 | 0.0 | 19600.0 | 18976.0 | 0.0 | 20124.0 | 0.0 | 15170.0 |

month

- January 1st is the least safe day of the year.
- The 1st day of every month is the least safe day of the month because this is usually the day when people are required to pay their rent and other recurring bills. These financial stressors often lead to commiting crimes.

**Time Series Visualizations:**

To model an ARIMA for the given data, it was imperative to visualize the time series of the crime count in NYC.
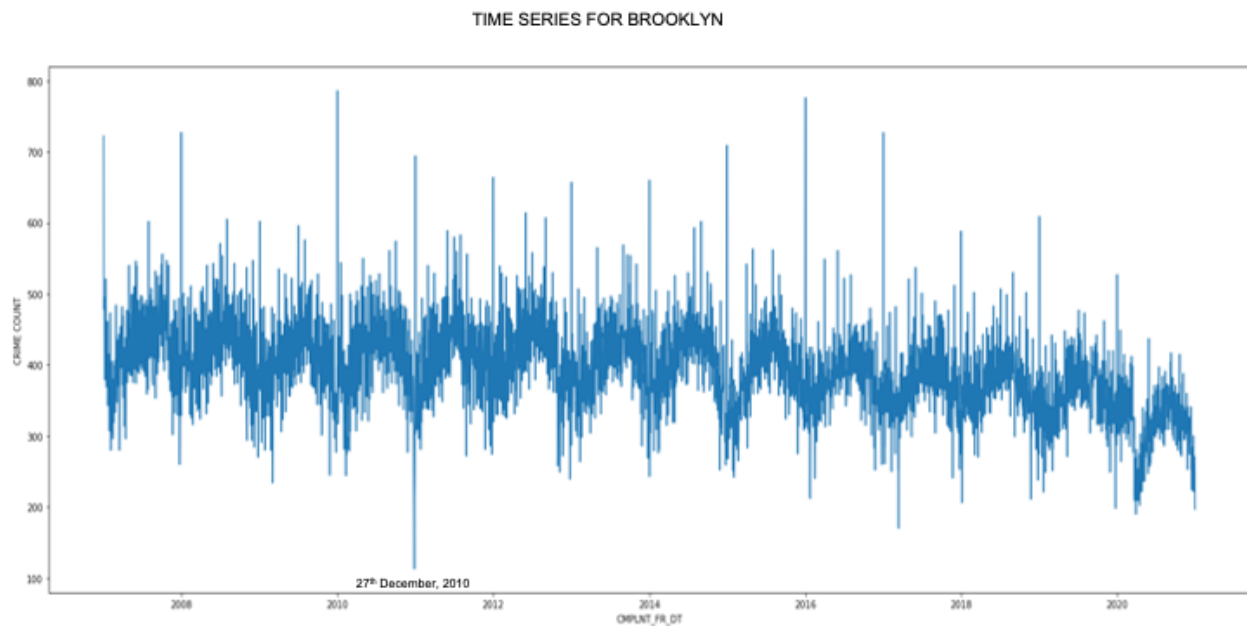


We also wanted to check whether similar patterns hold true when the data set is split borough-wise, and to analyze the trends in each borough.
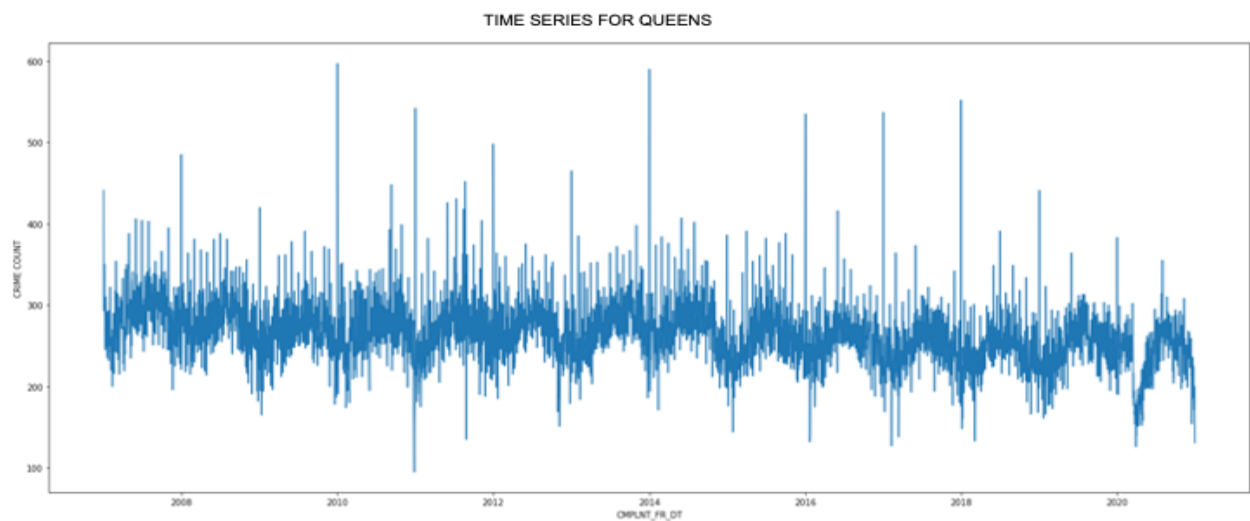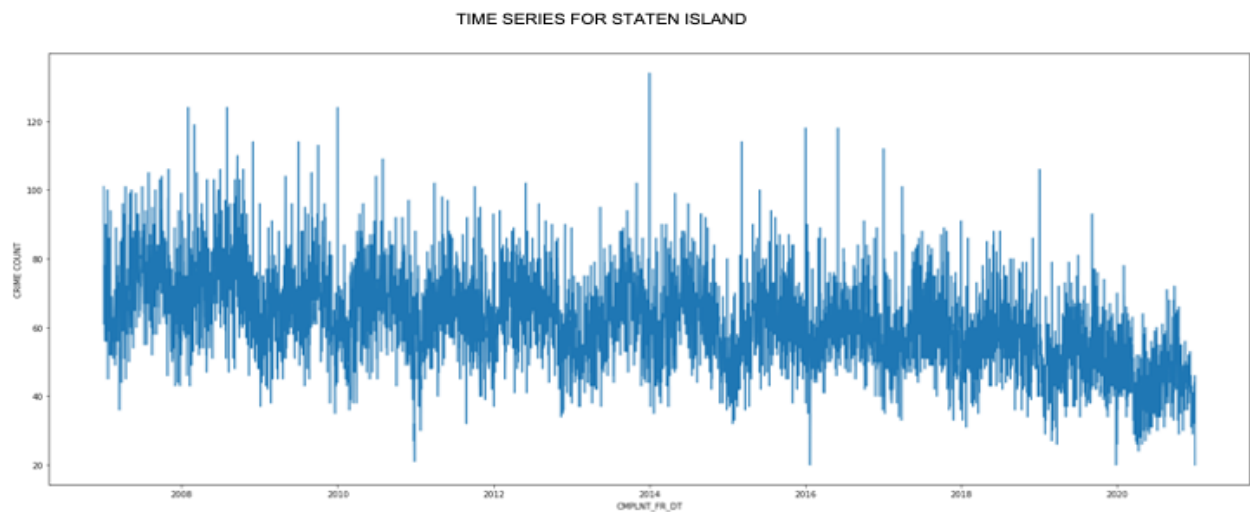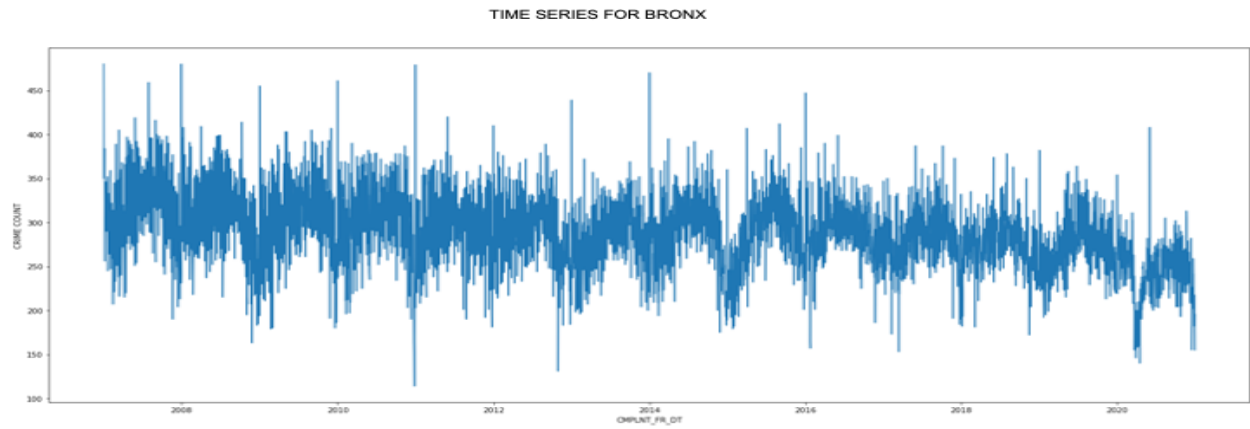


- The peak in the graph denotes 1st June, 2020. This was a result of the widespread protests and unrest caused by the death of George Floyd. Many shops were looted and some sites also experienced protest violence. The table below shows the crimes which were committed on 1st June, 2020.

| OFNS_DESC | COUNT |
|---|---|
| BURGLARY | 313 |
| CRIMINAL MISCHIEF & RELATED OF | 185 |
| POSSESSION OF STOLEN PROPERTY | 65 |
| PETIT LARCENY | 59 |

TIME SERIES FOR BROOKLYN



- The drop in crime count on 27th December, 2010 is common throughout all the boroughs. New York City witnessed one of its worst snow storms in its history on this day which reduced the mobility of the people.

TIME SERIES FOR BRONX



TIME SERIES FOR STATEN ISLAND



TIME SERIES FOR QUEENS

From the three time series above, it can be observed that the time series of crimes per day is very similar across all boroughs. There are a few outliers but in general the shape of the graph is similar.

**AutoRegression Integration Moving Average (ARIMA) Model**

ARIMA is a class of models that explains a time series of events based on its own past values. These values are used to forecast the future values. One of the main reasons we went with the ARIMA forecasting model is because of its robustness. The model requires a long historic horizon of data. Our dataset had complaints dating all the way back to 2006 which were over 6 million entries. This allows the algorithm to have a historic demand which helps in a more accurate forecasting.

The ARIMA is a combination of two different models - AutoRegression (AR) and Moving Average (MA) model. The blend of these two models do not overfit the data during training. The ARIMA model is characterized by 3 terms: p, d, q.

- p is the order of the AR term. The p term for the AR model can be determined by looking at the Partial Autocorrelation (PACF) plot since the PACF can be seen as a correlation between the series and its own lag.
- d is the number of differences required to make the time series stationary.
- q is the order of the MA term. It refers to the number of lagged forecast errors that should go into the model. Just like we determine the p value by analyzing the PACF plot, we can find the q value by looking at the Autocorrelation (ACF) plot.

There are two basic types of ARIMA models that are used for forecasting. They are Non-Seasonal ARIMA and Seasonal ARIMA (also known as SARIMA). The non-seasonal ARIMA consists of just the above three characteristics whereas the SARIMA is used with seasonal data.
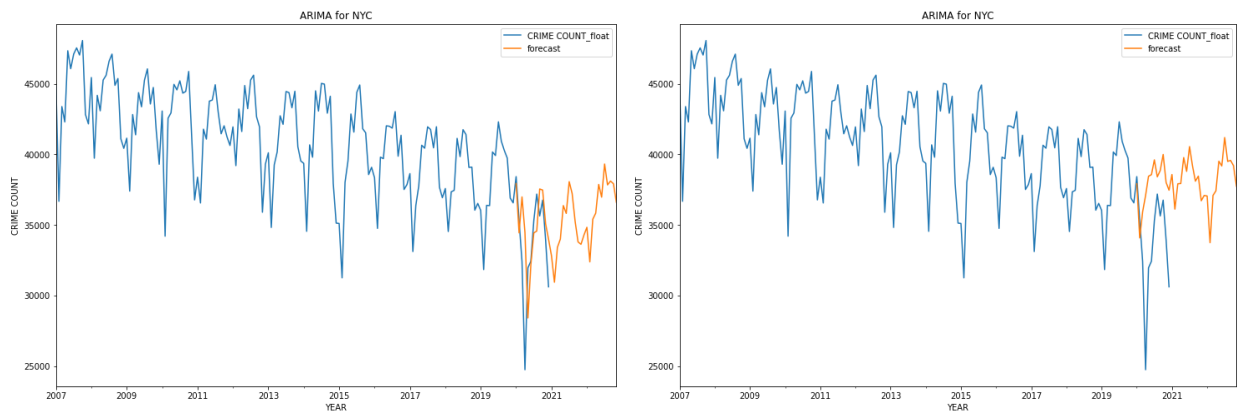
We add three new hyperparameters to the AR, I and MA components of the model as well as an additional parameter to determine the seasonality of the period of seasonality. We then configure the following four parameters:

- P: This is the seasonal autoregressive order of the model
- D: Seasonal difference order
- Q: Seasonal moving average order
- m: The number of time steps that we want to make for our model. We could select m to be 3 (for quarter of a year), 12(one year) or maybe more depending on how we want to analyze our time series.

Once these parameters and hyperparameters are determined, we start to fit our data into the model and train it for forecasting. Then the results for each borough were observed for the years 2021 and 2022.
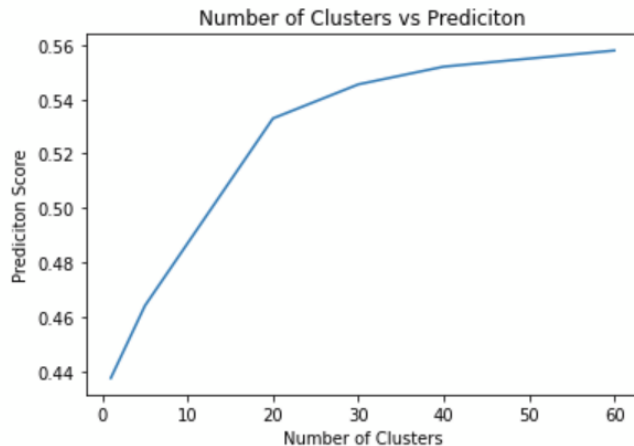
## Results

- Due to the spread of the COVID-19 virus, there was a significant drop in crime rate in the year 2020. This also led to a Mean Squared Error Percentage of 7.5%. If this pandemic had not occurred, the Mean Squared Error Percentage for this model would have been 2.5% giving us an overall accuracy of 97.5%.
- After a significant drop in crime numbers in 2020, we see that the number of crimes that are expected to take place in the next two years will increase in all five boroughs of New York City.



*Difference between forecasting due to presence and absence of COVID-19 respectively*

## Clustering with KNN

We decided to further our forecasting research by evaluating if we could categorize our future predictions. We have decided to use a KNN model to provide these classifications due to its low training time and classification accuracy given the robustness of our data set. With the number of features in the data set, we decided to do some data preprocessing and limit the number of features used in the KNN, as described in the Data Exploration section. After not finding anything significant in the Pearson Correlation Matrix, we tried running the KNN model on all of the features. This approach did not work given our limited computational power. To overcome this we started to remove features from our subsequent trials based on the cardinality of each feature. In our many tries to train a predictive KNN model, we noticed that the classification accuracy of our test data hovered around 20%. For our final model, the features we settled on using were those associated with the ARIMA data; time features and location features. Given all of this we trained our final data set over a variety of k clusters, giving us a graph of prediction percentages:

Number of Clusters vs Prediciton

Using the 60 clusters here gave us a predictions rate of 55%, which we deem is fairly accurate given our data set. Trying to use this model with our prediction from ARIMA led to our main pitfall.

## Pitfalls

ARIMA has proven to do a good job in forecasting crime trends in each of the boroughs for the next year, however our goal of using this information for preventive measures did not come to fruition. In order to use the KNN that we trained, we need to have specific instance predictions instead of projections. Our expectation was we would be able to get the data needed to predict not only the total number of crimes per month, but also have instance data that would fit into our KNN Model, giving us the predicted classification of crime for time and location. With advancements in modeling techniques for time series data, we expect to complete all of our goals for this project.

## Conclusion

Given the robustness of our dataset, we have successfully projected an accurate forecast of crime in New York City. Through our analysis, we verified that one of the most accurate predictors of future crime is the historical crime rates. We have seen that the outside factors, such as blizzards and protests, have significant effects on crime rate that are hard to predict. We furthered our analysis by producing a KNN Classification model that efficiently and somewhat confidently predicted crime classification given a certain amount of features. This work lays the groundwork for future research into this problem, starting with research into different time series analysis techniques, and if they would be able to predict instances rather than features. Papers have already been published on "Predictive Policing", leaving us with future routes of taking this work and furthering it for the safety of others. Given the success we have had so far, we believe that this is the start of great predictive models that can help with the prevention of crime.

**References:**

(1)  =  https://www.neighborhoodscout.com/ny/new-york/crime#data