

A PROJECT REPORT
On
“CRIME PREDICTION AND ANALYSIS”

Submitted to
ITM (SLS) Baroda University

Made By

B.Tech AI 5 th SEM	MUSTAFA LANEWALA	22C22033 (AI)
B.Tech AI 5 th SEM	MURTUZA MAHUDAWALA	22C22011 (AI)
B.Tech AI 5 th SEM	ABDEALI KHANRAHIM	22C22009 (AI)

UNDER THE GUIDANCE OF
PROF. Dr. PRADEEP LAXKAR



SCHOOL OF COMPUTER SCIENCE AND ENGINEERING(SOCSET)
ITM (SLS) BARODA UNIVERSITY

November-2024

ABSTRACT

We are interested in applying machine learning methods to datasets regarding crime (crime statistics in particular cities) and possible related factors (such as tweet data, income, etc.). Specifically, we are interested in investigating if it is possible to predict criminal events for a specific time and place in the future (for example, assigning a risk level for a shooting within the next week to different neighbourhoods).

To be better prepared to respond to criminal activity, it is important to understand patterns in crime. In our project, we analyse crime data from the Toronto Dataset, scraped from publicly available in Kaggle.

The use of AI/ML in predicting crimes or an individual's likelihood for committing a crime has promise but is still more of an unknown. The biggest challenge will probably be "proving" to politicians that it works. When a system is designed to stop something from happening, it is difficult to prove the negative. Companies that are directly involved in providing governments with AI tools to monitor areas or predict crime will likely benefit from a positive feedback loop. Improvements in crime prevention technology will likely spur increased total spending on this technology. We also attempt to make our classification task more meaningful by merging multiple classes into larger classes. Finally, we report and reflect on our results with different classifiers, and dwell on avenues for future work.

Keywords— Python; Machine Learning; Clustering; Time Series; Education; Students; Performance;

CONTENT

1. Introduction

- 1.1. Problem Statement
- 1.2. Scope And Objectives:
- 1.3. Proposed Model

2. Literature Review

- 2.1. Summary:

3. Basic Concepts/Technology Used

- 3.1. Machine Learning:
- 3.2. Classification:
- 3.3. Time Series Analysis

4. System Requirements

5. Data Dictionary

6. Implementation

- 6.1. Creating the core functionality

7. Verification Table

8. Conclusion and Future Work

9. References

Chapter 1

Introduction

Crimes are a common social problem affecting the quality of life and the economic growth of a society. It is considered as an essential factor that determines whether or not people move to a new city and what places should be avoided when they travel. With the increase of crimes, law enforcement agencies are continuing to demand advanced geographic information systems and new data mining approaches to improve crime analytics and better protect their communities. Although crimes could occur everywhere, it is common that criminals work on crime opportunities they face in most familiar areas for them. By providing a data mining approach to determine the most criminal hotspots and find the type, location and time of committed crimes, we hope to raise people's awareness regarding the dangerous locations in certain time periods. Therefore, our proposed solution can potentially help people stay away from the locations at a certain time of the day along with saving lives. On the other hand, police forces can use this solution to increase the level of crime prediction and prevention. For police resources allocation. It can help in the distribution of police at most likely crime places for any given time, to grant an efficient usage of police resources. By having all of this information available, we hope to make our community safer for the people living there and also for others who will travel there.

1.1. Problem Statement:

Our study aims to find spatial and temporal criminal hotspots and also forecasting of crime using a set of real-world datasets of crimes. We will try to locate the most likely crime locations and their frequent occurrence time. In addition, we will predict what type of crime might occur next in a specific location within a particular time. Finally, we intend to provide an analysis study by combining our findings of a particular crime's dataset with its demographic's information

1.2. Scope And Objectives

Our proposed solution can potentially help people stay away from the locations (crime hotspot) at a certain time of the day along with saving lives. On the other hand, police forces can use this solution to increase the level of crime prediction and prevention. For police resources allocation. It can help in the distribution of police at most likely crime places for any given time.

- Technology is Noticeable
- Elimination of Confusion
- TO help people be aware of the crimes and help the society
- Interactivity
- Crime forecasting

1.3. Proposed Model

In this work, we will build a machine learning module. The model works on the concept of Time Series Forecasting and Clustering. After successful running of the module the analysis and the forecasting results are shown through graphs and plots. We intend to provide an analysis study by combining our findings of a particular crimes dataset with its demographic's information.

Chapter 2

Literature Review

Machine Learning is undeniably one of the most influential and powerful technologies in today's world. More importantly, we are far from seeing its full potential. There's no doubt, it will continue to be making headlines for the foreseeable future. This article is designed as an introduction to the Machine Learning concepts, covering all the fundamental ideas without being too high level.

Machine learning is a tool for turning information into knowledge. In the past 50 years, there has been an explosion of data. This mass of data is useless unless we analyse it and find the patterns hidden within. Machine learning techniques are used to automatically find the valuable underlying patterns within complex data that we would otherwise struggle to discover. The hidden patterns and knowledge about a problem can be used to predict future events and perform all kinds of complex decision making.

Most of us are unaware that we already interact with Machine Learning every single day. Every time we Google something, listen to a song or even take a photo, Machine Learning is becoming part of the engine behind it, constantly learning and improving from every interaction. It's also behind world-changing advances like detecting cancer, creating new drugs and self-driving cars.

To learn the rules governing a phenomenon, machines have to go through a **learning process**, trying different rules and learning from how well they perform. Hence, why it's known as Machine Learning.

Ada Lovelace, one of the founders of computing, and perhaps the first computer programmer, realized that **anything in the world could be described with math**.

Around 200 years later, these fundamental ideas are critical in Machine Learning. No matter what the problem is, it's information can be plotted onto a graph as data points. Machine Learning then tries to find the mathematical patterns and relationships hidden within the original information.

More importantly, this meant a mathematical formula can be created to derive the relationship representing any phenomenon. Ada Lovelace realised that machines had the potential to understand the world without the need for human assistance.

2.1. Summary

The literature review begins by defining Machine Learning (ML) as a branch of artificial intelligence focused on developing algorithms that allow computers to learn and make decisions from data without explicit programming. It discusses how ML originated, tracing its roots back to early pattern recognition work and advancing rapidly in recent years due to increased computing power, larger datasets, and improved algorithms. ML's influence spans a wide range of fields, including healthcare, finance, retail, and entertainment, where it plays a crucial role in driving innovations and optimizing operations.

Chapter 3

Basic Concepts/Technology Used

Machine learning is an application of artificial intelligence (AI) that provides systems the ability to automatically learn and improve from experience without being explicitly programmed. **Machine learning focuses on the development of computer programs** that can access data and use it learn for themselves.

The process of learning begins with observations or data, such as examples, direct experience, or instruction, in order to look for patterns in data and make better decisions in the future based on the examples that we provide. **The primary aim is to allow the computers learn automatically** without human intervention or assistance and adjust actions accordingly.

Types of Machine Learning Algorithms:

There are 4 types of Machine Learning today:

1. Supervised Machine Learning Algorithms
2. Unsupervised Machine Learning Algorithms
3. Semi-Supervised Machine Learning Algorithms
4. Reinforcement Machine Learning Algorithms

1. Supervised machine learning algorithms can apply what has been learned in the past to new data using labelled examples to predict future events. Starting from the analysis of a known training dataset, the learning algorithm produces an inferred function to make predictions about the output values. The system is able to provide targets for any new input after sufficient training. The learning algorithm can also compare its output with the correct, intended output and find errors in order to modify the model accordingly.

2. In contrast, Unsupervised machine learning algorithms are used when the information used to train is neither classified nor labelled. Unsupervised learning studies how systems can infer a function to describe a hidden structure from unlabelled data. The system doesn't figure out the right output, but it explores the data and can draw inferences from datasets to describe hidden structures from unlabelled data.

3. Semi-supervised machine learning algorithms fall somewhere in between supervised and unsupervised learning, since they use both labelled and unlabelled data for training – typically, a small amount of labelled data and a large amount of unlabelled data. The systems that use this method are able to considerably improve learning accuracy. Usually, semi supervised learning is chosen when the acquired labelled data requires skilled and relevant resources in order to train it / learn from it. Otherwise, acquiring unlabelled data generally doesn't require additional resources.

4. Reinforcement machine learning algorithms is a learning method that interacts with its environment by producing actions and discovers errors or rewards. Trial and error search and delayed reward are the most relevant characteristics of reinforcement learning. This method allows machines and software agents to automatically determine the ideal behaviour within a specific context in order to maximize its performance. Simple reward feedback is required for the agent to learn which action is best; this is known as the reinforcement signal.

Classification

In the scope of Machine Learning, classification is an approach of supervised learning where the result set is to be catalogued as one of many existing result classes which is already trained.

Various classification algorithms include-:

Logistic Regression
Random Forest Classification
Support Vector Classification
Decision Tree Classification

Time Series

Time series analysis comprises methods for analysing time series data in order to extract meaningful statistics and other characteristics of the data. Time series forecasting is the use of a model to predict future values based on previously observed values. While regression analysis is often employed in such a way as to test theories that the current values of one or more independent time series affect the current value of another time series, this type of analysis of time series is not called "time series analysis", which focuses on comparing values of a single time series or multiple dependent time series at different points in time. Interrupted time series analysis is the analysis of interventions on a single time series.

Different classical time series forecasting methods; they are:

2. Autoregression (AR)
3. Moving Average (MA)
4. Autoregressive Moving Average (ARMA)
5. Autoregressive Integrated Moving Average (ARIMA)
6. Seasonal Autoregressive Integrated Moving-Average (SARIMA)
7. Seasonal Autoregressive Integrated Moving-Average with Exogenous Regressors (SARIMAX)
8. Vector Autoregression (VAR)
9. Vector Autoregression Moving-Average (VARMA)
10. Vector Autoregression Moving-Average with Exogenous Regressors (VARMAX)
11. Simple Exponential Smoothing (SES)
12. Holt Winter's Exponential Smoothing (HWES)

Chapter 4

System Requirements

Jupyter Notebook Version

Version 7.2.0 is used in our Project.

The Jupyter notebook is a tool which we can use for our machine learning project and statistical analysis. We can download anaconda from the web source and within it Jupyter notebook most useful tool for machine learning purpose.

Python Version

Python 3.12.5 is used for this project.

Python is a very useful programming language. It is object oriented and interpreted. It is a high-level language. There are lots of in-built libraries in Python for machine learning purpose which we can use easily.

Windows Version

Jupyter notebook and python 3 can be used in all the operating systems including Windows, iOS and Linux.

It is best useful in Linux but can be used in windows as well.

It can be run on windows xp, vista, 7, 8 and the latest version windows 10 as well windows 11

Chapter 5

DATA DICTIONARY

Sl. no	Feature Name	Feature Description
1	Index	It stores the serial number of the crimes which are reported.
2	event_unique_id	Stores the unique id of the crime
3	occurrencedate	Date in which crime had actually occurred.
4	reportdate	Date in which the complaint was lodged.
5	premisetype	Type of the premise in which the crime took place like Commercial apartment, house, etc.
6	ucr_code	It is the abbreviation of Uniform Crime Reporting which enumerates offense codes. A code list that describes a criminal offense within a code book.
7	ucr_ext	It is a code list which describes the type of offense within a code book.
8	offence	Stores the type of offence such as Assault, robbery, etc.
9	reportedmonth	Month in which the complaint was lodged.
10	reportedday	Day in which the complaint was lodged
11	reporteddayofyear	Stores the report day of the year such as 99, 321, etc.
12	reporteddayofweek	In which day of the week the crime was reported.
13	reportedhour	Stores the hour in which the report was filed.
14	occurrenceyear	Year in which the crime took place.
15	occurrencemonth	Month in which crime took place
16	occurrenceday	Date of the month in which the crime took place.

17	occurrencedayofyear	Day of the year in which the crime took place.
18	occurrencedayofweek	Day of the week in which the crime took place.
19	occurrencehour	The hour in which the crime took place.
20	MCI	It is an abbreviation of Major Crime Indicators which store the type of crime that has been reported such as assault, break and enter, robbery, etc.
21	Division	Stores the division number of the neighbourhood.
22	Hood_ID	Contains the ID of the neighbourhood.
23	Neighbourhood	Name of the neighbourhood where the crime took place.
24	Long	Longitude of the place where crime had occurred.
25	Lat	Latitude of the place where crime had occurred.
26	ObjectId	It is an index of the crime after the reporting to the police.

Chapter 6

Implementation

Implementation is the process of converting the designed system architecture into working modules where it is made sure that all the functional and non-functional requirements are met.

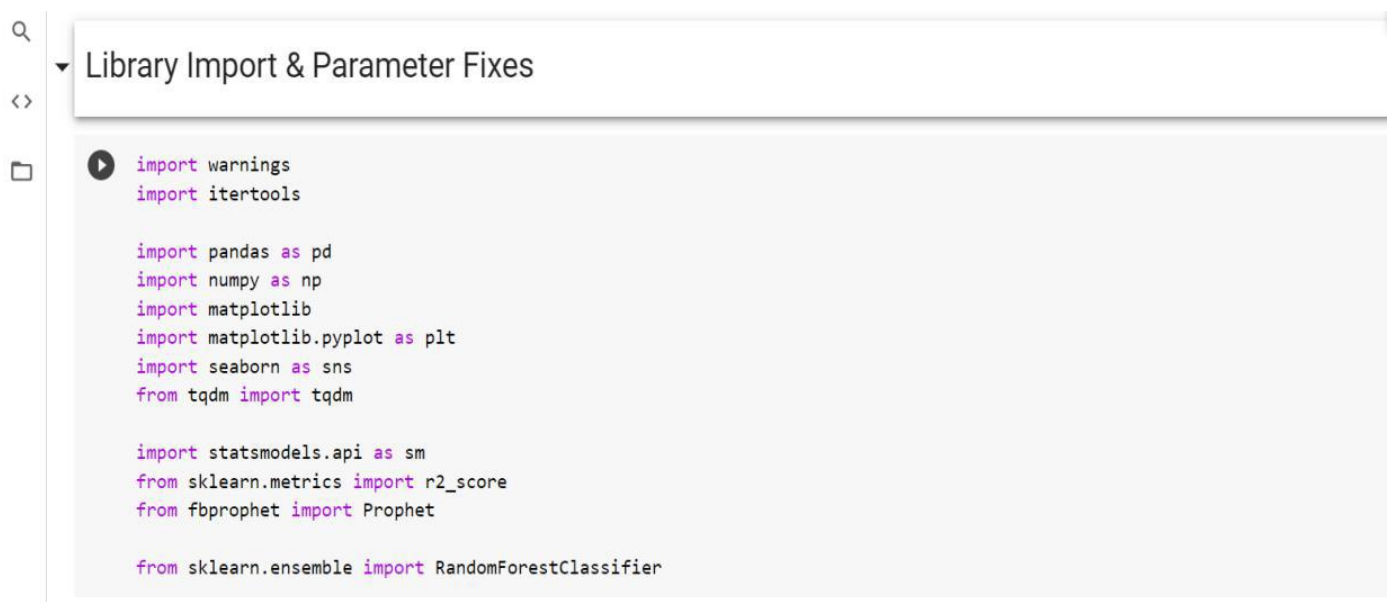
The implementation section is divided into two parts-

- Creating the Core functionality (Machine Learning Module)

7.1. Creating the core functionality

Dataset Importing

Importing of dataset from .csv file into the jupyter notebook.



The screenshot shows a Jupyter Notebook interface. On the left, there is a sidebar with a search icon, a code icon, and a folder icon. The main area displays a code cell titled "Library Import & Parameter Fixes". The code cell contains the following Python code:

```
import warnings
import itertools

import pandas as pd
import numpy as np
import matplotlib
import matplotlib.pyplot as plt
import seaborn as sns
from tqdm import tqdm

import statsmodels.api as sm
from sklearn.metrics import r2_score
from fbprophet import Prophet

from sklearn.ensemble import RandomForestClassifier
```

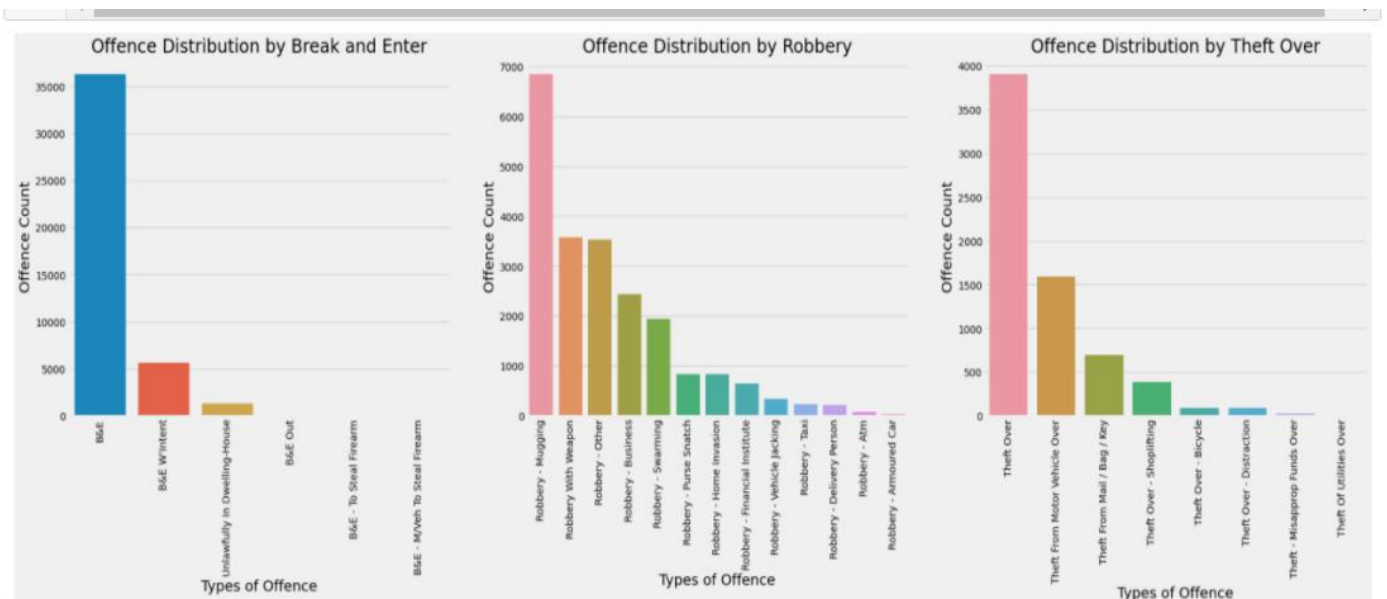
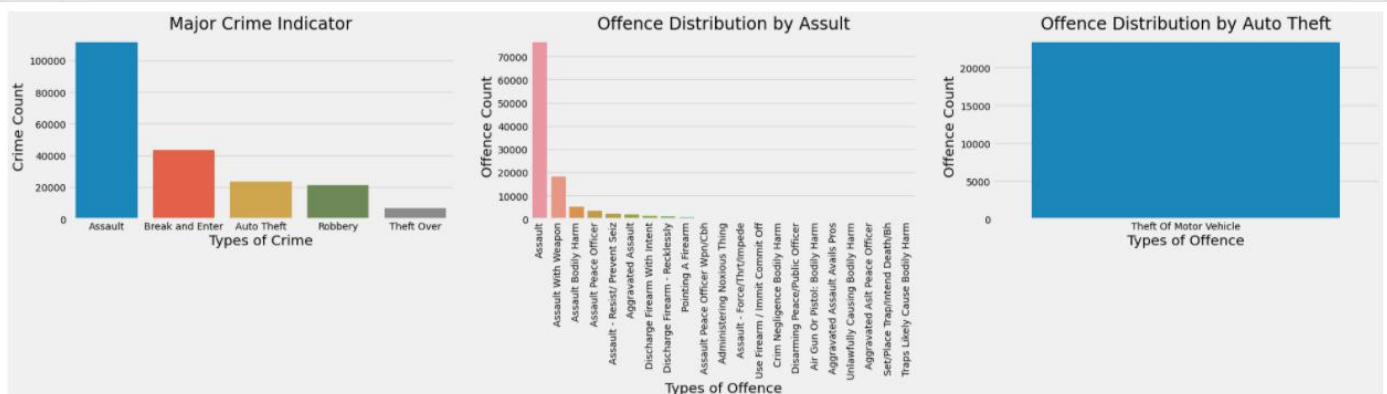
Our team members are well versed with data manipulation. The data analysis was done using Python3 in jupyter notebook. Some data that NULL in the dataset, they were dropped. The Data was thoroughly analyzed to get the the minimum error in our Time Series method. Encoding were also done. Then the clustering was done.

```
df=pd.read_csv("MCI_2014_to_2019.csv")
df['Total'] = 1
df.head()
```

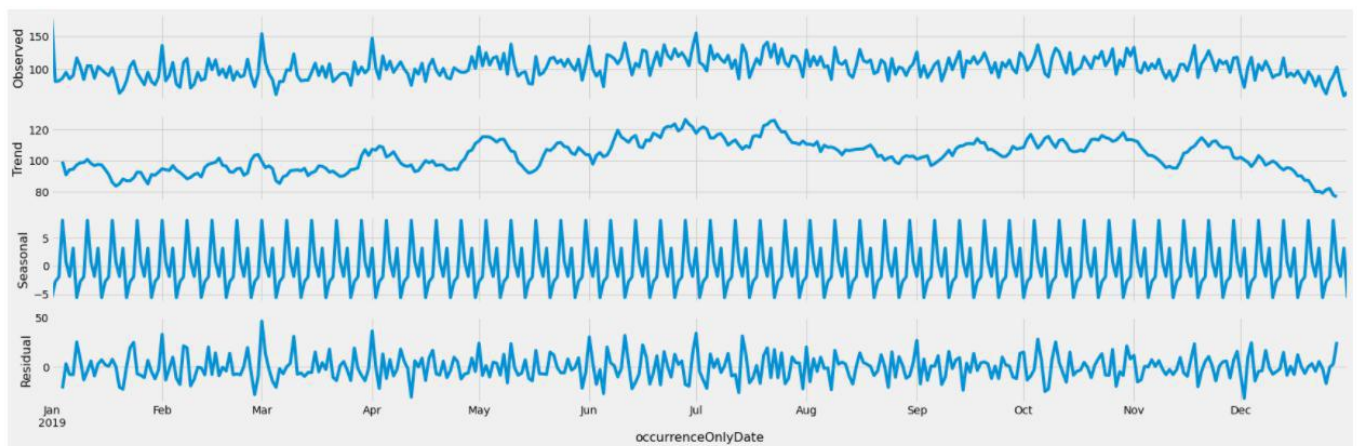
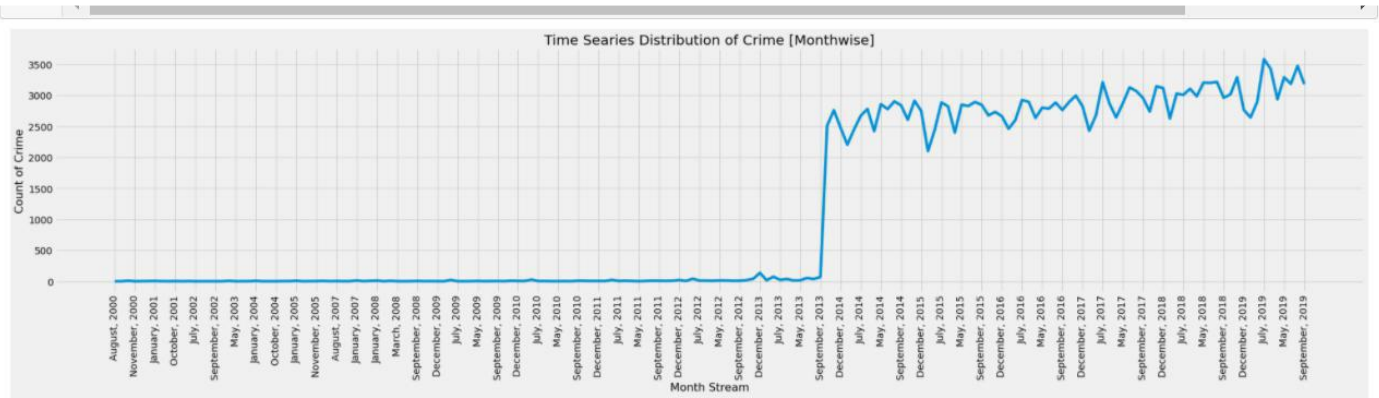
	X	Y	Index_	event_unique_id	occurrence date	reported date	premise type	ucr_code	ucr_ext	offence	...	occurrence day of week	occurrence hour
0	-79.405228	43.656982	7801	GO-20152165447	2015-12-18T03:58:00.000Z	2015-12-18T03:59:00.000Z	Commercial	1430	100	Assault	...	Friday	3
1	-79.307907	43.778732	7802	GO-20151417245	2015-08-15T21:45:00.000Z	2015-08-17T22:11:00.000Z	Commercial	1430	100	Assault	...	Saturday	21
2	-79.225029	43.765942	7803	GO-20151421107	2015-08-16T16:00:00.000Z	2015-08-18T14:40:00.000Z	Apartment	2120	200	B&E	...	Sunday	16
3	-79.140823	43.778648	7804	GO-20152167714	2015-11-26T13:00:00.000Z	2015-12-18T13:38:00.000Z	Other	2120	200	B&E	...	Thursday	13
4	-79.288361	43.691235	7805	GO-20152169954	2015-12-18T19:50:00.000Z	2015-12-18T19:55:00.000Z	Commercial	1430	100	Assault	...	Friday	19

5 rows × 30 columns

Plotting and analysis of different crime:



Time Series Analysis for Total Crime Count:



ARIMA Time Series Forecasting:

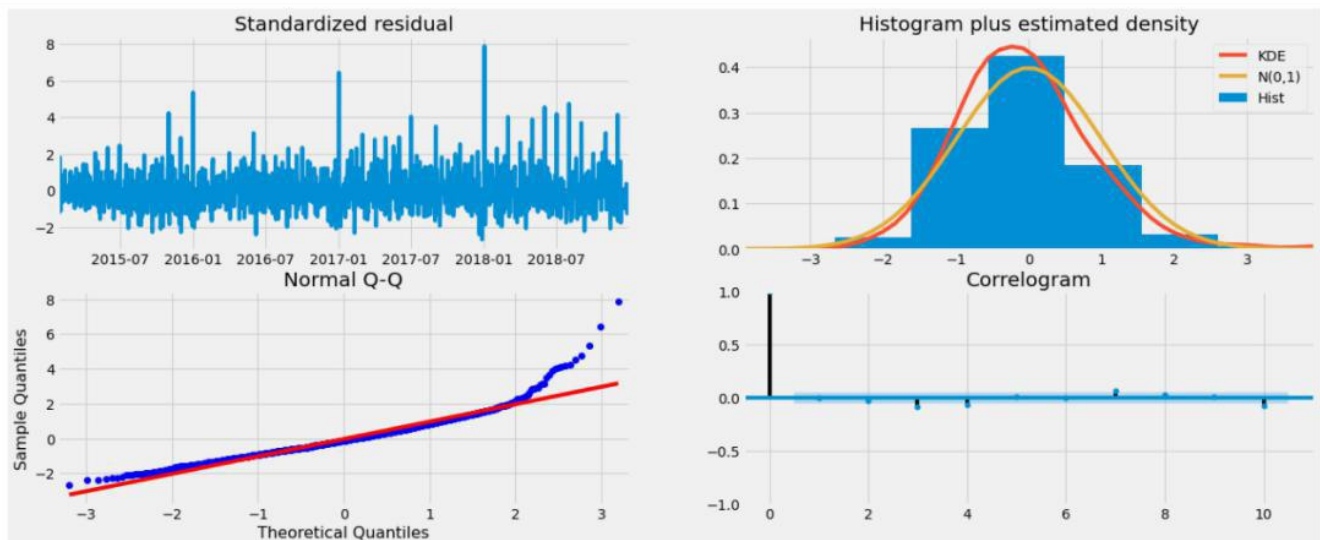
A popular and widely used statistical method for time series forecasting is the ARIMA model. ARIMA is an acronym that stands for Autoregressive Integrated Moving Average. It is a class of model that captures a suite of different standard temporal structures in time series data.

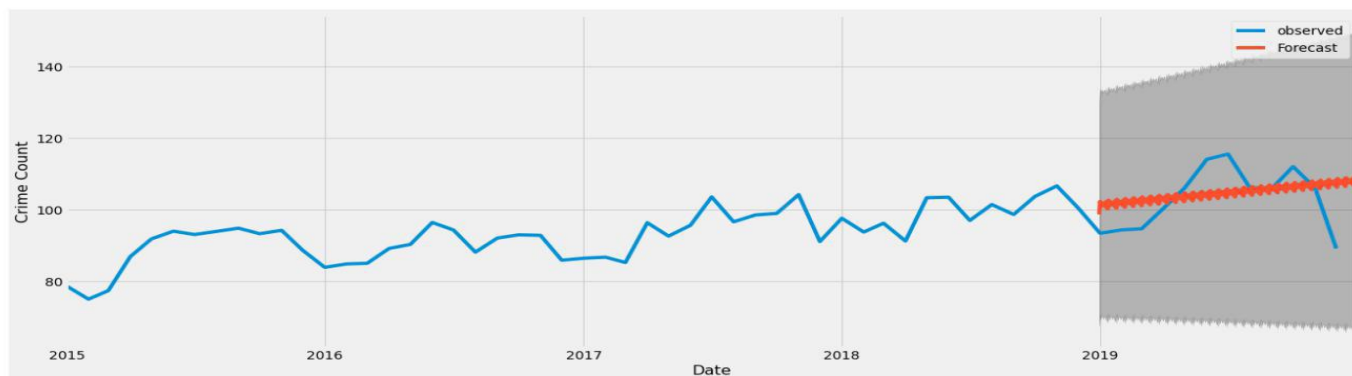
SARIMA for Time Series Forecasting:

An extension to ARIMA that supports the direct modelling of the seasonal component of the series is called SARIMA.

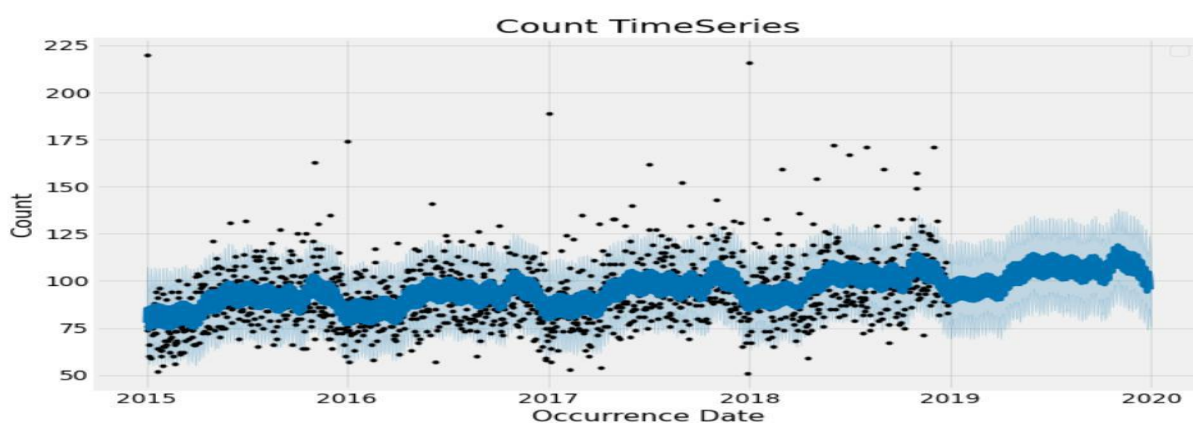
In this tutorial, you will discover the Seasonal Autoregressive Integrated Moving Average, or SARIMA, method for time series forecasting with univariate data containing trends and seasonality.

```
1 SARIMAXresults.plot_diagnostics(figsize=(20, 8))  
2 plt.show()
```



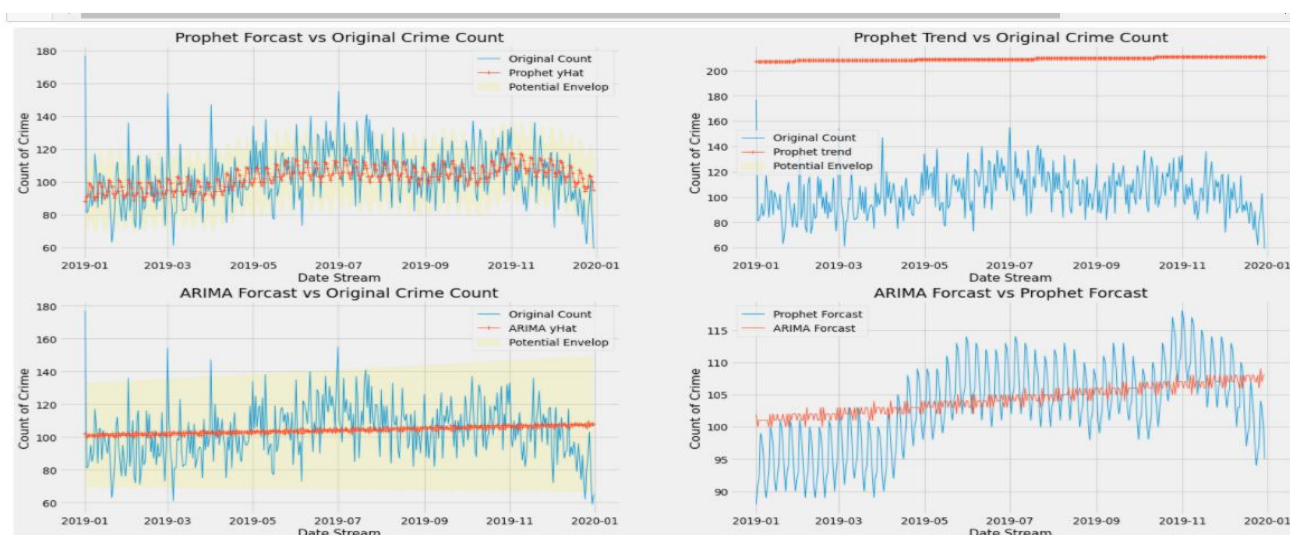


SARIMAX FORECASTING



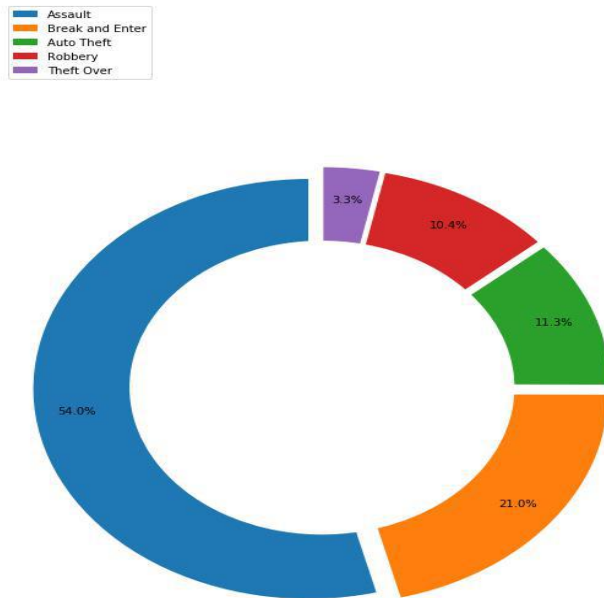
Model Comparison and Actual Value Deviation Analysis

Prophet Time Series Forecasting

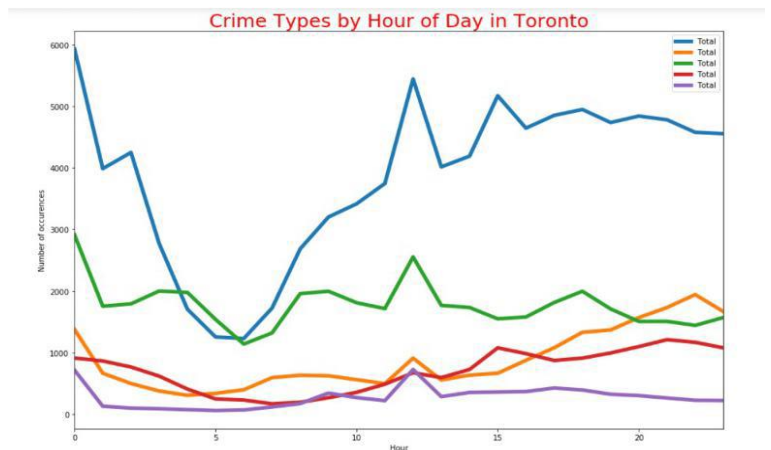
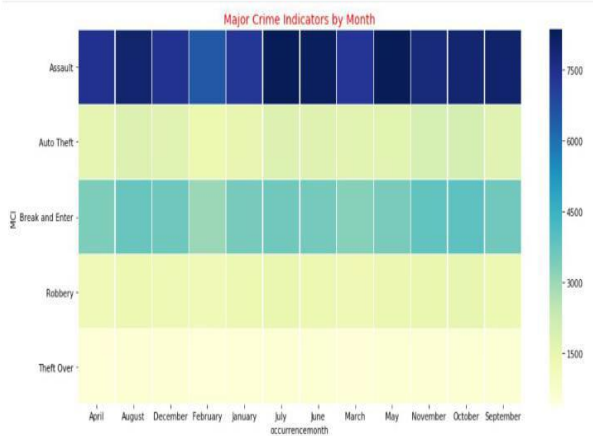
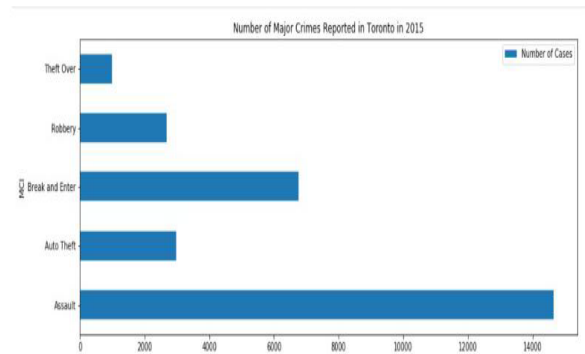
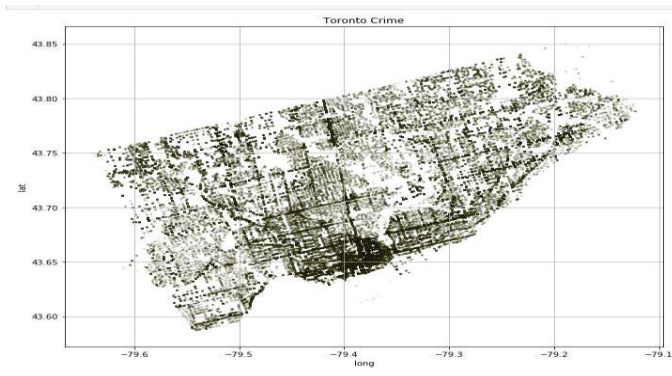


Differences between the original and the predictive model plotted in graphical manner

Share of Criminal Offences in Toronto



Different types of Criminal Offences in Toronto

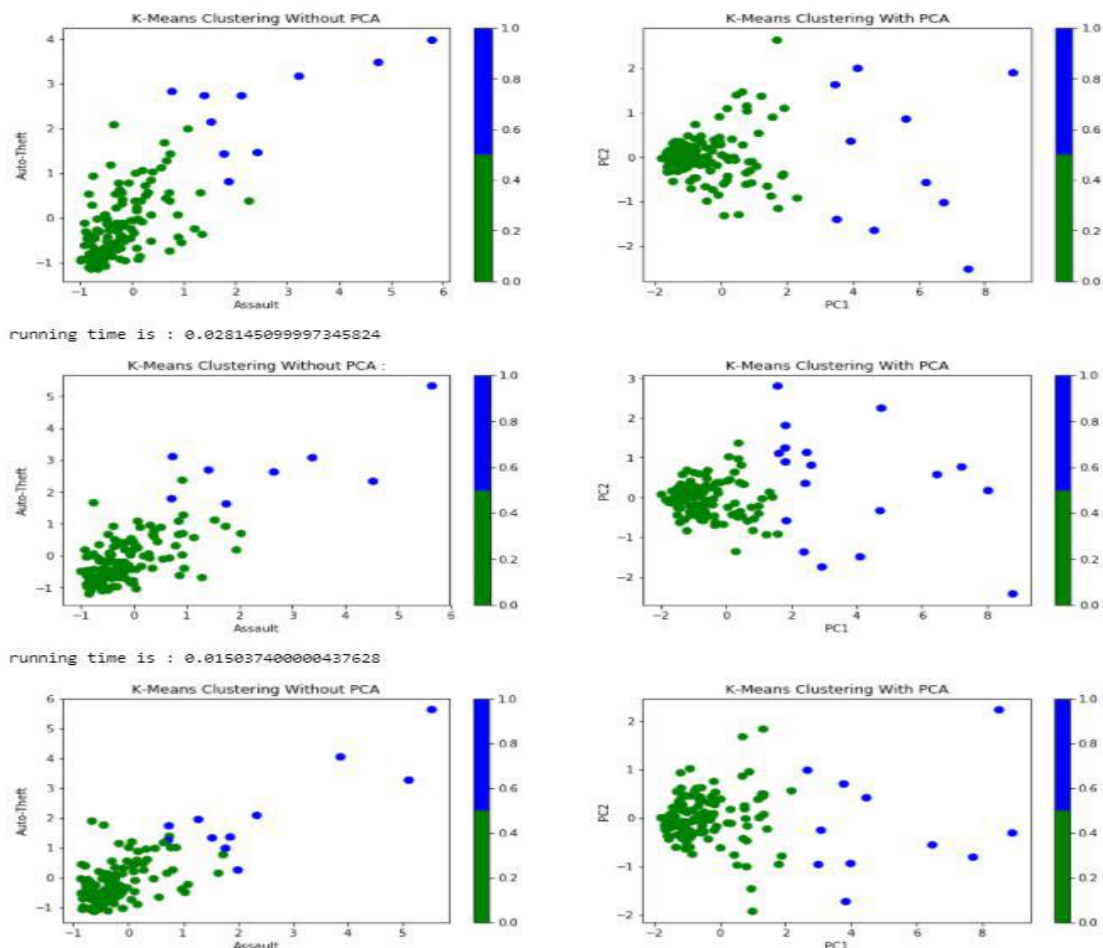


K-means Clustering:

Kmeans algorithm is an iterative algorithm that tries to partition the dataset into K predefined distinct non-overlapping subgroups (clusters) where each data point belongs to only one group. It tries to make the intra-cluster data points as similar as possible while also keeping the clusters as different (far) as possible. It assigns data points to a cluster such that the sum of the squared distance between the data points and the cluster's centroid (arithmetic mean of all the data points that belong to that cluster) is at the minimum. The less variation we have within clusters, the more homogeneous (similar) the data points are within the same cluster.

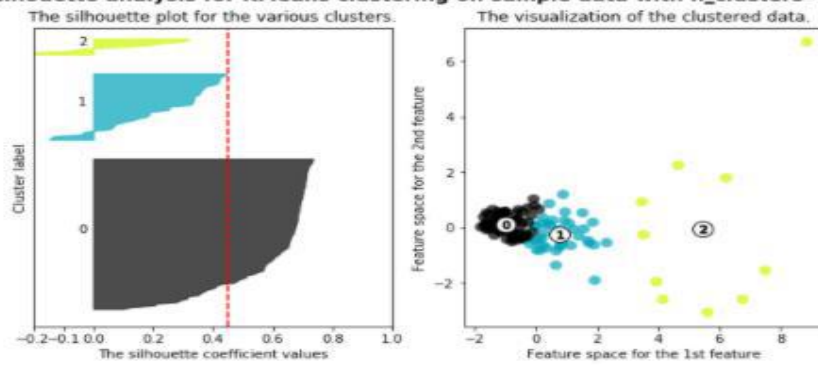
PCA:

PCA is used in exploratory data analysis and for making predictive models. It is commonly used for dimensionality reduction by projecting each data point onto only the first few principal components to obtain lower-dimensional data while preserving as much of the data's variation as possible.

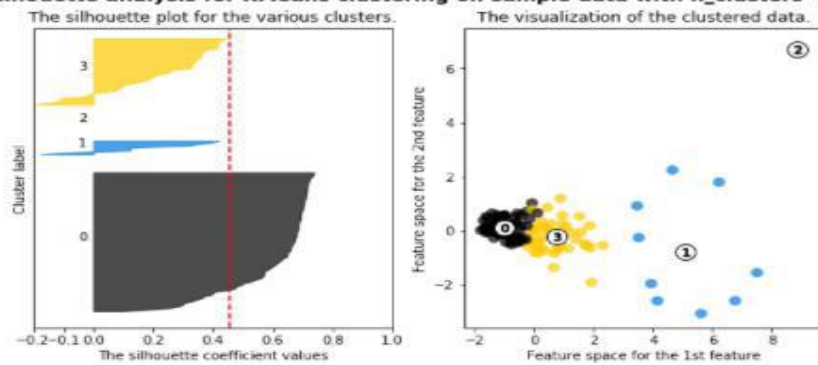


We have done clustering one without PCA and one with PCA and have found that applying PCA Before doing clustering can help in getting better clusters and also visualization becomes much better.

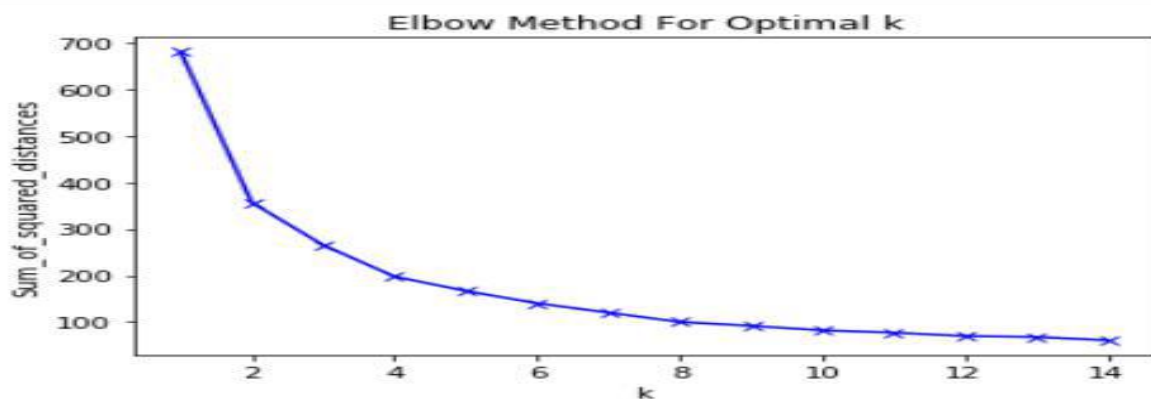
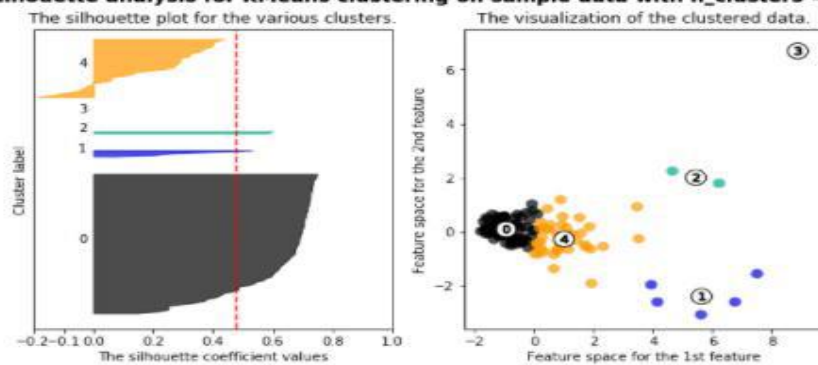
Silhouette analysis for KMeans clustering on sample data with n_clusters = 3



Silhouette analysis for KMeans clustering on sample data with n_clusters = 4



Silhouette analysis for KMeans clustering on sample data with n_clusters = 5



Elbow method was used to determine the number of clusters which should be used in order to get better clusters. It consists of plotting the explained variation as a function of the number of clusters. In this case we have considered 2 clusters.

Chapter 7

Verification Table

1	SARIMAX		Prophet	
2	MSE	RMSE	MSE	RMSE
3	290.11	17.03	248.27	15.76

So as we applied SARIMA model Prophet model and got this much Mean Square Error and Root Mean Square Error. We have got less error in the Prophet time series model. So, we can use this model for further works.

```
print("Accuracy of Random Forest with OneHotEncoder : ",accuracy_score(y_test, y_pred))
print(confusion_matrix(y_test_OH, y_pred_OH))
print(classification_report(y_test_OH,y_pred_OH, target_names=definition_list_MCI))
```

Accuracy of Random Forest with OneHotEncoder : 0.5838382934955169

[[19536	2298	228	21	861]
[4883	5581	27	14	211]
[3280	309	260	7	371]
[1164	327	20	17	118]
[2862	634	142	9	1543]]
	precision	recall	f1-score	support
Assault	0.62	0.85	0.71	22944
Break and Enter	0.61	0.52	0.56	10716
Robbery	0.38	0.06	0.11	4227
Theft Over	0.25	0.01	0.02	1646
Auto Theft	0.50	0.30	0.37	5190
micro avg	0.60	0.60	0.60	44723
macro avg	0.47	0.35	0.35	44723
weighted avg	0.57	0.60	0.56	44723

For n_clusters = 2 The average silhouette_score is : 0.7307306081429717
For n_clusters = 3 The average silhouette_score is : 0.4511660188781014
For n_clusters = 4 The average silhouette_score is : 0.4529368049556418
For n_clusters = 5 The average silhouette_score is : 0.4776823069308707

These are the accuracy of Random Forest Classification and silhouette scores of Kmeans clustering respectively. We have good better score in K-means clustering using 2 clusters.

Chapter 8

Conclusion and Future Work

While, there is little reason to believe that the crime rate will increase dramatically in the first decade of the 21st Century, given the anticipated increases in the globalization, sophistication, and organization of crime, one may conclude that the impact of crime on Western societies may be more severe than the one witnessed under a similar rate of crime in the past. The goal of any society shouldn't be to just catch criminals but to prevent crimes from happening in the first place.

1. Predicting future crime spots.
2. Predicting who will commit the crime.

Chapter 9

References

[1] N.Cristianini and J.Shawe-Taylor (2000). *An Introduction to Support Vector Machines*.Cambridge University Press, London.

[2] <https://towardsdatascience.com/machine-learning/home>

[3] Prophet Time Series : <https://machinelearningmastery.com/time-seriesforecasting-with-prophet-in-python/>

[4] <https://www.ibm.com/industries/government/public-safety/crimeprediction-Prevention>