

Analysis and prediction of crime in New York City

Chebaane Zeineb

Higher School of Communication of Tunis
zeineb.chebaane@supcom.tn

Arfa Sirine

Higher School of Communication of Tunis
sirine.arfa@supcom.tn

Abbes Malek

Higher School of Communication of Tunis
malek.abbes@supcom.tn

Arij Flihi

Higher School of Communication of Tunis
arij.flihi@supcom.tn

Abstract

Crime is a serious and pervasive social issue that exists everywhere. The rate of crime has significantly increased in recent years. Advanced systems and fresh ideas are required for enhancing crime analytics in order to protect communities in response to this rise. Even though effective real-time crime prediction lowers crime rates, it is still a challenging topic for scientists because crime incidences depend on numerous complicated elements. In order to anticipate crimes in New York City, a variety of visualization approaches and machine learning algorithms are used in this study.

A raw dataset was processed in the first step, and several visualization techniques were used to better comprehend the data and the relationships between the various variables. After that, a machine learning system was employed to forecast different crime categories depending on user input and location of the users. The final stage is to create a user interface with folium and flask to simplify user interaction. This github repository contains the finished code: <https://github.com/SirineArfa/New-York-City-Crime-Prediction-Project>

Keywords: *Crime Analysis; Crime prediction; Data Visualization; Crime Maps; Machine Learning; Classification; Folium; Flask; Python; Random Forest Classifier*

1 Introduction

Proper analysis of past crime data aids in crime prediction and subsequently assists efforts to lower the crime rate. Investigating crime reports and immediately spotting fresh patterns, series, and trends are all part of the analysis process. This study aids in the quick preparation of statistics, queries, and maps. Because criminals are active and tend to operate in their comfort zones, it is possible to forecast the type of crime they will do next because, if they succeed in committing the first crime, they are likely to repeat it.

The following offense is typically attempted in a similar place and at a similar time. Studies indicate that there is a significant likelihood of repeat, even though this may not be true in all situations, making crimes predictable. This project suggests creating a web application and visual interface for a crime prediction tool in Python utilizing a variety of libraries, including Flask for the user interface, Folium for interactive leaflet maps, Pandas for data processing, etc. The suggested framework employs multiple visualization techniques to demonstrate the trend in crimes and various machine learning algorithmic approaches to forecast crimes.

Data preprocessing, data visualization, and model construction are the most crucial phases, and they are covered in more detail in the following sections. In a nutshell, the preprocessing stage entails data cleansing and transformation. Finally, in the model-building phase, we employed the Random Forest Classification algorithm to categorize the crimes

that may occur in a specific location. The visualization phase creates numerous reports and maps for the diagnostic and analysis process.

2 Related Work

Crimes are a common social problem affecting the quality of life and the economic growth of a society [1]. It is considered an essential factor that determines whether or not people move to a new city and what places should be avoided when they travel [2]. Today, a high number of crimes are causing a lot of problems in many different countries. In fact, scientists are spending time studying crime and criminal behaviors in order to understand the characteristics of crime and to discover crime patterns. Dealing with crime data is very challenging as the size of crime data grows very fast, so it can cause storage and analysis problems. In particular, issues arise as to how to choose accurate techniques for analyzing data due to the inconsistency and inadequacy of these kinds of data. These issues motivate scientists to conduct research on these kinds of data to enhance crime data analysis. The objective of this research is to apply suitable machine learning algorithm on crime data to predict the likelihood of a county having low, medium or high violent crimes

2.1 Crime analysis

Criminology is an area that focuses on the scientific study of crime and criminal behavior and law enforcement and is a process that aims to identify crime characteristics.[3] It is one of the most important fields where the application of data mining techniques can produce important results. Crime analysis, a part of criminology, is a task that includes exploring and detecting crimes and their relationships with criminals. The high volume of crime datasets and also the complexity of relationships between these kinds of data have made criminology an appropriate field for applying data mining techniques. Identifying crime characteristics is the first step for developing further analysis. The knowledge that is gained from data mining approaches is a very useful tool which can help and support police forces. The proposed framework provides visualization techniques that consider the location and many other information introduced by the user to predict the type of crime to overcome these limitations

2.2 Why Crime is Predictible

There is a strong body of evidence to support the theory that crime is predictable (in the statistical sense) mainly because criminals tend to operate in their comfort zone . That is, they tend to commit the type of crimes that they have committed successfully in the past, generally close to the same time and location. Although this is not universally true, it occurs with sufficient frequency to make these methods work reasonably well. There are major theories of criminal behavior, such as routine activity theory, rational choice theory, and crime pattern theory. These theories are consolidated into what is referred to as a blended theory.

A previous work that inspired us to look more into crime prediction is crime prediction based on weather, crime data, and temporal data [4]. In the paper, the authors employed feature selection techniques to determine the most significant features mainly the most occurred crimes and the correlation between the features, in forecasting crime calculations and rates in New York City over 5 years. They used both machine learning and deep learning techniques and provided benchmarking based on the prediction accuracy. Another interesting work that motivated us is spatiotemporal crime forecasting using Amsterdam police Data [5] in which they focused on Crime history variables, Environmental variables, Demographic variables, Socio-economic variables, and Proximity variables to provide more detailed and reasonable comprehension and prediction that highlights the reasons of the committed crimes.

3 Methodology

In this section, we explain our methodology on how to build a machine learning model and cross-validate it on the New York crime data. We want to predict the different categories of crime based on : time , victim description and location, thus, we implement this workflow :

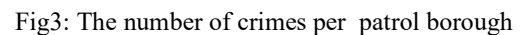
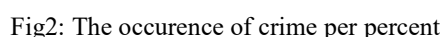
- Dataset extraction : This work relies on NYPD Complaint Data Historic dataset. This dataset includes different categories of crimes reported to the New York City Police Department (NYPD) from 2006 to 2019. The dataset contains 6901167 complaint and 35 columns including spatial and temporal information about crime occurrences along with their description and penal classification.
- Data Preparation : To further understand the data in hand and analyze the different distributions and relations between features, we undergo an Exploratory data analysis in order to answer questions about what, where and when crimes occur, and dealing with null values, outliers and unnecessary characteristics in order to clean the data and get the best possible accuracy for our machine learning model.
- Feature extraction: To select significant features we used the documentation provided with the data to select the features that we are going to need in our work then we used those features to create additional features such as correlation matrix , encoding techniques and detailing the time and dates like year , month , day and time zone.
- Modeling : After the data pre-processing, in order to classify three different types of crimes : felonies, misdemeanors and violations, based on their severity, we applied Random Forest model, The classification is mainly used to recognize the labelled classes by knowing their attributes in the dataset, thus predicting the class label for instances with known features. Hence, using the classifiers in crime prediction constructs a future-oriented model to identify the criminal type within a specific time.
- Model evaluation : we used the confusion matrix and the

Result visualization : we created a dashboard that presents the crimes prediction distributed in a map .

After studying the description of the features provided by our NYPD dataset we decided to keep just 12 features which are: CmplntFrDt, CmplntFrtm, OfnsDesc, LawCatCd, XCoordCd, YCoordCd, Latitude, Longitude, PatrolBoro, VicageGroup, VicRace and VicSex. Those features present the victim description date time location and the type of the crime which are the necessary information that we need to predict the type of the crime based on location and the person characteristics. We opted to fill missing values based on the distribution of the values in the data-set. As for the timestamp values we replaced all the nulls with the median value in each column and deleted outliers with unreasonable values in years and ages.

Fig1: NYPD dataframe after cleaning

For this part, we built derived values from our initial data which are more informative and non redundant. We started with generating year, month and day columns based on CMLPLNTRDT. Then from CMLPLNTRTM we categorized the different daytime into four classes: morning, afternoon, evening and night. The same way, having 14 types of crimes we grouped them into only 15 classes thus our prediction will be processed according to a fewer number of classes.



OneHotEncoding : All of PATROLBORO, VIC- SEX and VICRACE were encoded using this tech-nique as these variables don't present any natural order to take into consideration.

3.3 Random Forest Classification Model:

which is a meta estimator that fits a number of de- cision tree classifiers on various subsamples of the dataset and uses averaging to improve the predic- tive accuracy and control over-fitting and we used sklearn python library to implement the model. After training, we evaluate this model which gave us : 0.519 as score of the accuracy.

0	100506	9923	40640	8	1490	4770	27052	28260	11442	9572	262	12203	131	1575	59
1	14348	72200	0	4	839	16673	6917	179	8770	9130	112	22207	194	1203	5
2	35021	0	18790	9	842	423	14349	21336	4091	6799	153	1821	0	3080	15
3	2	0	0	167903	0	0	0	0	370	29	1	0	0	0	0
4	932	455	434	1	15736	67	199	256	185	201	9	285	8	24	0
5	2441	6515	33	4	48	10619	2083	251	2388	2405	1420	3831	32	1098	4
6	11569	1358	16187	1	196	1422	24944	6629	2462	2983	1308	1429	8	998	8
7	10269	17	8426	1	202	140	5553	11504	2043	860	77	136	0	609	4
8	3411	2848	462	97	108	1798	1484	599	6816	876	593	1710	35	367	1
9	1854	2513	459	2	71	1789	1557	243	862	5915	911	1871	13	258	15
10	496	35	123	11	36	2420	1817	367	1470	1176	10791	399	2	1317	27
11	9696	10329	109	3	400	4348	2127	113	2219	3603	97	19492	94	248	21
12	6	5	0	0	0	2	1	0	7	1	0	6	9	0	0
13	180	91	171	3	14	766	302	160	478	89	581	126	2	2292	0
14	0	0	1	0	0	0	1	0	1	4	6	0	0	2	17
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14

Fig4: Random Forest Confusion matrix

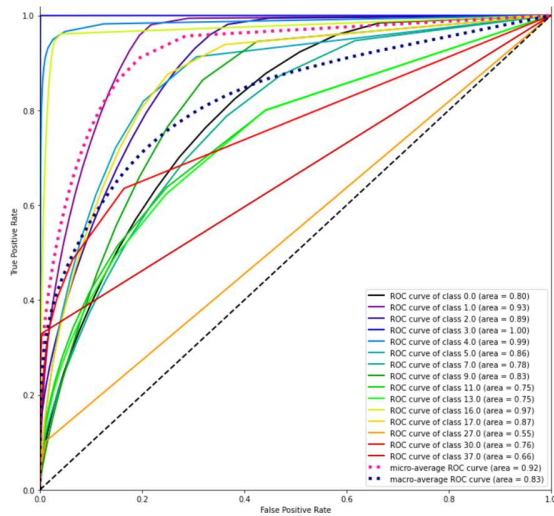


Fig5: Random Forest ROC curve

3.4 User Interface:

we built a web application using Flask and Folium to allow the user to interact with the map and predict the type of crime that could happen. The user can enter his gender, race, age, the date and hour in which he wants to predict the type of crime, the location on the map and finally, the place. This information is then transformed to fit the model input, and then, using the loaded model weights file, we predict the type of crime and send it back to the user along with the potential subtypes of that crime.

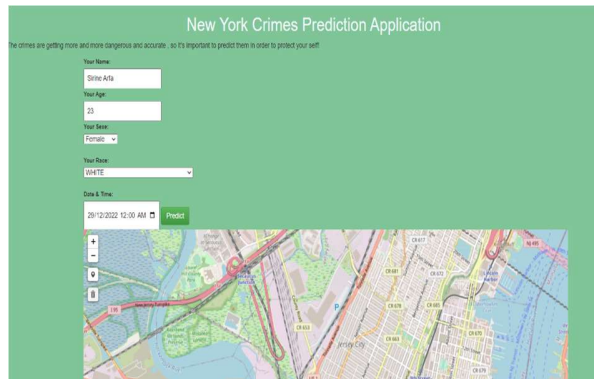


Fig6: User Interface

Attention: you might encounter the following crime with the highest probability:



Crime Type	Probability
aggravated assault	0.3056666666666667
grand larceny	0.06666666666666667
petty larceny	0.07272727272727273
burglary	0.10647117144444444
robbery	0.05748181818181818
sexual assault	0.0000000000000000
public safety crimes	0.0000000000000000
adultery	0.0000000000000000
drugs and alcohol crimes	0.0000000000000000
death and robbery	0.0000000000000000
kidnapping	0.0000000000000000
terrorism	0.0000000000000000
child abuse crimes	0.0000000000000000

Fig7: Crime Prediction

4 Conclusion

The perception of a community as crime ridden can deter people from going there and induce residents to move away. This causes damage to the economy. Crimes affects the economy by placing a financial burden on taxpayers and governments because of increased needs for police, courts and corrections facilities, as well as intangible costs including psychological trauma and reduced quality of life for crime victims.

Hence, many researchers tried to solve it and predict the most criminal hot spots to increase the understanding of dangerous places at certain times.

In this work, we used the data-set provided by the NYPD. As a first step we explored the data in order to understand its pattern and produce insights. Then We kept the most essential features by performing techniques of feature selection and feature extraction. After that, we applied our Machine Learning Model Random Forest. As a result Random Forest was very close in prediction accuracy for 0.55. Finally, We have created a user interface to enable users to enter their information and get the class of the crime that can happen in a particular location at a specific time.

5 References

- [1] Bogomolov, B. Lepri, J. Staiano, N. Oliver, F. Pianesi and A. Pentland, 'Once Upon a Crime, Towards Crime Prediction from Demographics and Mobile Data', CoRR, vol. 14092983, 2014.
- [2] R. Arulanandam, B. Savarimuthu and M. Purvis, 'Extracting Crime Information from Online Newspaper Articles', in Proceedings of the Second Australasian Web Conference - Volume 155, Auckland, New Zealand, 2014, pp. 31-38.
- [3] Malathi A., Santhosh B.S., Algorithmic Crime Prediction Model Based on the Analysis of Crime Clusters; Global Journal of Computer Science and Technology; Volume 11 Issue 11 Version 1.0 July 2011.

[4] Elluri, Lavanya Mandalapu, Varun Roy, Nirmalya. (2019). Developing Machine Learning Based Predictive Models for Smart Policing. 10.1109/SMARTCOMP.2019.00053

[5] Rummens, Anneleen Hardyns, Wim Pauwels, Lieven. (2017). The use of predictive analysis in spatiotemporal crime forecasting: Building and testing a model in an urban context. Applied Geography. 86. 10.1016/j.apgeog.2017.06.011.

[6]