

CRIME ANALYSIS AND PREDICTION
SYSTEM PROJECT USING
MACHINE LEARNING TO SOLVE CURRENT
TRENDS OF VIOLENT CRIMES

BSCLMR230620

October 4, 2021

Contents

1	INTRODUCTION	3
1.1	Background	3
1.2	Problem Statement	4
1.3	Objectives	4
1.3.1	General Objectives	4
1.3.2	Specific objectives	4
1.4	Justification	5
1.5	Scope	5
1.6	Limitation	5
2	2 Literature Review	6
2.1	Theoretical Review	6
2.2	Similar Projects	7
2.2.1	Crime Prediction System using K-means clustering	7
2.2.2	Crime Prediction System using Linear regression	9
2.2.3	Crime Prediction System using Kernel Density Estimation method	9
2.3	Conceptual Framework	10
3	Methodology	12
3.1	Data Source	12
3.2	Data collection tools	12
3.2.1	Reports	12
3.2.2	Surveys	12
3.3	System design	12
3.3.1	Use case diagram	12
3.3.2	Architectural Framework	13
3.4	Implementation	14
3.4.1	Administrator	14
3.4.2	User Interface	14
3.4.3	Registration module	15
3.4.4	Dataset	15
3.4.5	Prediction algorithm	15
3.4.6	Crime Database	15
3.5	Relationship between conceptual framework and architectural framework	15
3.6	Tools	16
3.6.1	Python	16
3.6.2	Flask framework	16
3.6.3	Scikit library	16
3.6.4	Numpy	16
3.6.5	Pandas	16
3.6.6	Mysql	16
3.7	Testing	16

3.7.1	Administrator testing	16
3.7.2	Prediction algorithm testing	17
3.7.3	Crime database testing	17

1 INTRODUCTION

1.1 Background

The concept of crime as examined as a reason for predictive policing which was aimed to stop crime before it occurred. Crimes committed back then became widespread since there were no prediction tools that could be used to automatically identify where the crimes were committed. Not even trackers were developed and this mortified people in the neighbourhood. Technology that was first developed as an evolution of computers (super computers) which would simply store data but not a large amount of data. [16] There were murder sprees [2] in the neighbourhood, burglaries, homicide and rape crimes and yet the hotspots that these crimes were committed were not identified. [3].

These crimes would occur again and as for such reasons, later in the early 1900s there were developments of the models which would show where the crimes occurred. A major disadvantage to this was brought about by the computers that were used; they could not input big data becoming difficult to handle. Ada Lovelace [4], a mathematician described a sequence of operations and Charles Babbage punch-card machine and used equations to develop the patterns of crime to solve such activity that have become a thoughtful way to simulate the activity. In 1950, Allan Turing proposed a hypothesis imposing that machines can be used in convincing humans if a machine can succeed and achieve artificial intelligence. [5] Bernard Wodrow and Marcian Holf [6] developed neuron models that would be used to detect the binary patterns and this would be used in world application.

Later in 1967 the nearest neighbour algorithm was written which would allow the computer to use pattern recognition. The K-NN algorithm was introduced by Fix and Hodges in 1951 which was used for classification and regression [7]. In the 1990s, work on machine learning was shifted from knowledge based to a data driven approach and scientists later developed programs for computers to analyse large amount of data. Over the decades there has been many improvements of this algorithm there are new approaches that have emerged making it more improved and efficient. Later on in the current technology using this algorithm, alongside ordinal regression, logistic regression and decision stamp are implemented.

A crime predictive software developed was the CrimeScan software to predict violent crime [8]. The first trial of this software to test its effectiveness was when the researchers tackled violence and identify the individuals within the crime hotspots.

For example my case study in Atlanta could be in possible causes. If there was a robbery in the neighbourhood, a high probability was that they would use non-lethal weapons such as a baseball bat but later a similar crime would happen they would use a gun. Many cases would arise and as a result mathematical models would be used to predict crime and research has shown that the equations once developed, they can pinpoint the crime hotspots. Later on research done by Neill and Gorr [9] to develop a software tool for crime prediction which

would use the former dataset(historical dataset) which would show where the violent crimes would occur. [10]

As technology has evolved so has the artificial intelligence. How we face crime today using machine learning has made it easier to identify the criminal hotspots and the patrol to act right away. From the beginning of the fuzzy K-NN algorithm to the modern and more developed K-NN algorithm that uses other algorithms such as additive regression and linear regression algorithm that have been tested thus improving efficiency. [11] Machine learning as it has recently trend with its improvement of other algorithms such as logistic regression that has helped us understand dependent and independent variables have become tremendous tool for crime pattern detection whereby these crime patterns are automatically identified. Also we have the PredPol software was created to show the scientific analysis with the goal of data that could spot criminal behaviour and show crimes that may likely occur during a particular period. [1]

1.2 Problem Statement

The main problem addressed in this project is to discover the crime hotspots in the area using the dataset. Through documentation cases and investigation discovery channels, machine learning and data science has made it easier to predict crime and make the work easier. At times, many cases go cold and unsolved, while others have remained a mystery for even us as humans to contemplate the occurrence of events especially in a crime event. People in the society would want answers to know what really happened at the scene of the crime, parents want justice for their children and vice versa. This has further led to the criminals run free because of insufficient evidence and later on for such cases they remain unsolved. For example if a crime has occurred in an area that involves both murder and burglary, it can be assumed that maybe two suspects are involved and this makes have separate cases while it turns out that one suspect has committed both of the crimes.

1.3 Objectives

1.3.1 General Objectives

The aim of this project is to utilise and make the crime prediction features present in the main dataset extracted from the main link and run it with the help of the machine algorithm.

1.3.2 Specific objectives

The main objective is to train a model for prediction which will be trained and tested using the dataset. The enhancement and improvement of the model will be done using the K-NN algorithm(K-Nearest-Neighbour) and visualization for better accuracy. Visualisation of the dataset is done to analyse the crimes done while the algorithm will be used for the crime prediction. These are the procedures that will help achieve the main goal. First is data collection and

in this case I use the dataset provided for a certain area. The data will indicate occurrence of events (time, day and week). Secondly is to classify which will rely on the solved crimes. I will then implement the pattern identification which will show the trends in crime and the patterns as well. The prediction will use a decision tree to test on the attributes and will help the algorithm make a better decision. Lastly, the visualization whereby the crime hotspot area is represented geographically using the heat map to show the levels of activity based on the color images and analyse the data I want.

1.4 Justification

This project will benefit the society and even the state at large since criminal activities will be controlled. Handling and solving of these cases will make the society feel safe and make the socio-economic lifestyle easier and safer. Not only will it benefit the state but even the patrol making it easier to know the criminal events taking place at a specific time and act fast. Data included in this case will be used to find the unknown patterns from known data and facts and the approach will be automated by learning of the crime series.

1.5 Scope

The scope of the project is to improve the efficiency and accuracy of the machine learning algorithms in the prediction of the violent crimes and identifying the crime patterns when determining crime hotspots and learning the criminal trends. Law enforcement agencies can be warned and prevent crime from occurring by implementing more surveillance within the prediction zone. Furthermore, with complete automation, it overcomes the drawbacks of the current system and the law enforcement can depend on more of these techniques even in the future. Designing the machine learning algorithm to identify patterns of such crimes can be the start of the future of this study. Despite the systems having a big influence on the current crime prevention, with the start of this project, this could be the next revolutionary change in crime rate.

1.6 Limitation

The project will face the following limitations: analysis of data which may be difficult due to inconsistency and incomplete, limitation to getting crime data records from the law enforcement which they may be reluctant to give detailed data and lastly the accuracy of the program which mainly depends if the training set is accurate. In addition to this, finding the patterns and trends can be a challenging factor since it takes a lot of crime, scanning it to see if it fits in a particular pattern. Some of these trends and patterns are difficult in an unspecified geographical area. To expound on this, there are areas whereby the system may not predict since it may not be available or marked on the heat maps.

2 Literature Review

2.1 Theoretical Review

The first related machine learning work was done in 1950 whereby the related work used mainly symbolic data and algorithm that was based on logic done by Alan Turing [5] and many of these crimes have been obtained by the authorities or even a survey. The data collected would be transferred time and location as well as the crime committed to the police departments and the crime data using the websites. The information later on would take a specific time period to analyse the

Tools used to predict time and location were violated and even sharing the datasets could vary. Grouping of the datasets back then to understand criminal behaviour and analysis took overtime because there were no developed algorithms back then. Trying to verify or even know where the crime was committed was difficult. At some point the coverage area was too large [13] to even know the ground cover and machine learning by then was under-developed.

The traditional method to detect crime hotspot in the area from the past data in the distribution of crime cases and as for others to assume the past pattern to be repeated in the future. A common method used was the KDE method which used to identify stable hotspot areas. This method was used based on the on trends to outperform the method and found that the forest algorithm was more efficient than the traditional KDE. To analyse and understand the crime pattern techniques there were data mining techniques and tools to understand knowledge from them. One of the methods or technique to discover the crime patterns. [14] High number of crimes over the years have forced the government to use technology to predict and control the crime. That was then the start of using machine learning algorithm to detect these crime patterns.

The first computer learning program and it was improved whereby it simulates the thought of the human brain. The Nearest Neighbour allowed the computer to identify pattern recognition. Gerald Dejong (1981) explained the concept on the based learning whereby the computer analysis trained data and the general rule disposing unimportant data. The machine Learning shifts from knowledge driven to the data approach that the computer programs analyse data.

The nearest neighbour algorithm developed at the beginning was developed as an algorithm for mapping of routes before it later advanced before being advanced in crime. Machine learning models are created and have prompted a new array of concepts. As used in the study, the machine learning repository have used different classes of algorithms that are used to identify the dataset which include the regression, classification and clustering. [15]

K-NN algorithm uses the similarity feature to predict the values of any data whereby the value assigned is based on how they closely related hence the average values are taken to be the final prediction and forecast in the crime area. There are techniques used which is the supervised learning technique and assumes the similarity of cases and aligns them into the category which are

similar or available. As a classification technique it evaluates the value of the function hence in the varsity of different scales which normalizes the training of data thus improving the data accuracy. [13] There are steps when conducting the crime analysis:

The system design of the proposed approach using the data mining. Data scientist and analyst have related to machine learning ordinary by use of the K-NN classifier to enhance its exactness with the relative computational time. The K-NN rule has been the most of the most regulated calculation. It holds the preparation set during learning and its easiest type. The key component of the algorithm is that they have similar qualities if grouped in a similar classification. The architecture explains the crime examination framework, structure and test the new framework. [16]

It shows the result of the algorithm of crime prediction system consisting of a collection of thick crime areas and a set of related crime forecaster. This operate primarily for the large region where the large amounts of individuals live and demonstrates that the suggested strategy achieves excellent precision over rolling time horizons in spatial and temporal crime forecasting. This paper's working process collects raw data the hotspot uses after splitting the data to create the new hotspot model and shows the predictive crime rate.

The Neural Network is used for the prediction's precision. The precision rate by using Neural Network the model accuracy is 60, 93 and 97. It presents a geographical analysis-based and self-regressive approach to automatically identify large danger urban crime areas and to represents crime patterns in each region reliably. Show the result of the algorithm of crime prediction system consisting of a collection of thick crime areas and a set of related crime forecaster. [14]

2.2 Similar Projects

These are some of the projects that have been implemented in this field

2.2.1 Crime Prediction System using K-means clustering

The crime analysed and considered homicide taking into account and the trend has been descending from 1990-2011 and the clustering technique extracts information in form of the crime data setting using the Rapid Miner Tool since its a solid and complete package with the flexible support option and the figure below shows the proposed architecture and flow diagram. [17].

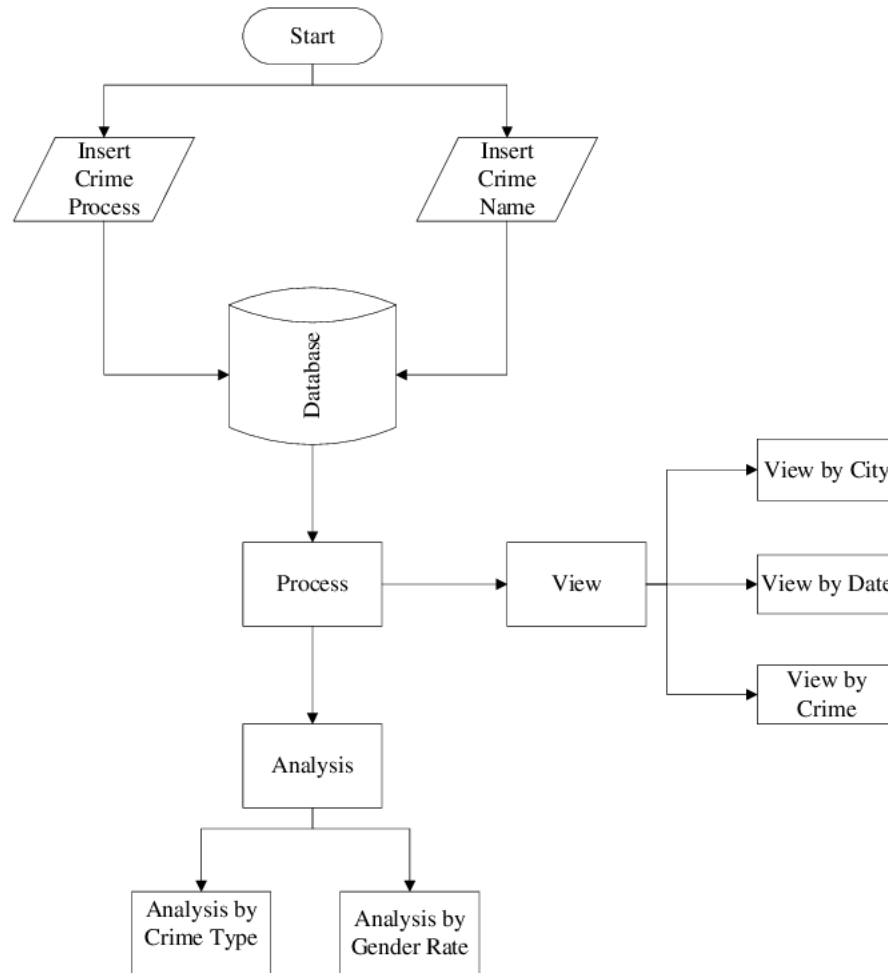


Figure 1: K-means clustering proposed architecture

The system in this case can predict regions with high probability of crime occurrence and can visualize the crime prone areas. The modules explain this process: Start The system starts : we input the crime process and the crime name which takes the data and information to the database. The system inserts the crime process and the crime name to the database. Database The database section :in this case consist of the data retrieved from the sources and at the process stage its preprocessed to make sure that the data is improved and later viewed. View Once viewed: the data is categorised by the place(city) when it occurred (date), specific duration it happened(time) and the type of crime. Analysis After the viewing of the data the collected data is now analysed into two: By the crime type whereby crime is categorised into different types for example rape, robbery, kidnapping and homicide and also by the gender rate(male or female) The machine algorithm will now train the dataset.

2.2.2 Crime Prediction System using Linear regression

Using two separate categories of the visualisation tools whereby one provides the geographical view of crimes and the visualisation ability of the social networks. [?]. Giving description of crime which includes the location and personal characteristics as input. By using this technique it involves some of the integrated components:geographical profiling, social network analysis,crime patterns and physical matching. Figure 2 shows the system design of the Linear regression technique.

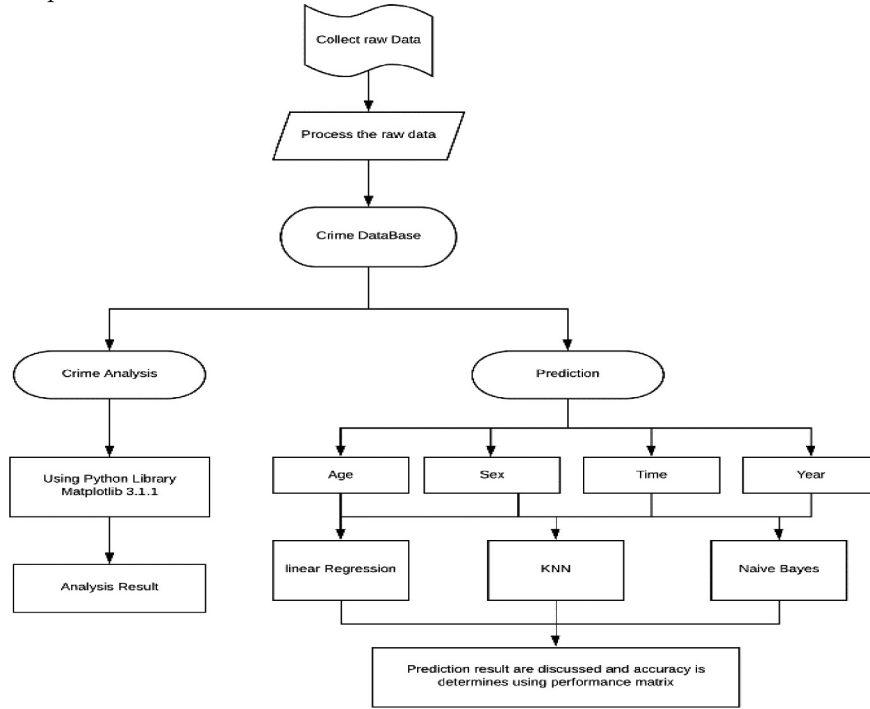


Figure 2: System Architecture: Linear Regression

From the flow diagram above data collected by use of reports and surveys is processed and stored in the crime database. Data is then classified whereby it involve crime analysis and uses python library and if use by machine learning prediction, the dataset used to train consist of the age, sex(gender), time when the crime happened and the year as well. the data is now trained using different algorithms and prediction models and give the prediction result.

2.2.3 Crime Prediction System using Kernel Density Estimation method

This focuses on the criminal hotspot of criminal events in a geographical space which might be related to the criminal theories and provide basic mechanism for the events and the police to use information. It may include criminal theories and are considered as the theoretical basis of crime prevention. Crime pattern

theories mostly explain the distribution of criminal events through the KDE method, a network approach that increases prediction accuracy. It improves the high-dimensional data and extracts the characteristics of data. In this case figure 3 shows a crime map from 2003-2013 with the crime rate of crime hotspot location. [19]

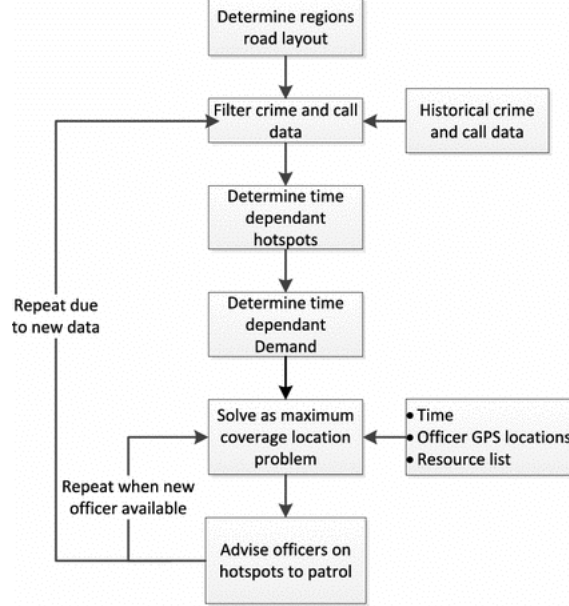


Figure3:KDE graph spatial mapping technique

Initially we identify the crime hotspot and this is achieved by taking historical crime data, lettering it in order to discard data not relevant to the problem, and performing kernel density estimation. Secondly, we use these hotspots as possible locations to send officers who are in patrol and identify which hotspots are chosen and what is determined by finding the configuration with maximum coverage of possible demand, using historical call data to predict the demand. The method adopted to calculate the optimal configuration of hotspots to patrol a version of the method Performing this analysis once however does not give a long-term solution to the patrolling problem as it is a dynamic problem. The location of hotspots and demand are time dependant, as is officer availability. Also once a hotspot has been patrolled the effect of deterring crime in the area has an unpredicted lifetime and hence the area may need to be visited again

2.3 Conceptual Framework

The figure below shows a conceptual framework for the crime prediction system. It will consist of the following: The data set that we will train our data. Secondly, the raw data will be improved at the data preprocessing which then the data train test will split. In this case, there are different techniques that can be used which include; Linear regression, K-NN algorithm, random forest

classifier and the neural network. Once data is ready, it will be trained using machine learning models and the data will be classified.

Figure 4: Conceptual framework diagram

The flow begins extracting data from the data collection on different roles of the data repository and the primary data is preprocessed and processed to criminal data. Data is preprocessed to make sure enhancement and improve experimental data. Once at the crime database the crime is classified to different models by the KNN classification. The crime is classified and identified by K-NN classification and crime prediction classified in Kernel and weighed K-NN.

3 Methodology

3.1 Data Source

The project will use the dataset obtained from Atlanta criminal records and included; classification of different types of crime in Atlanta; robbery, domestic violence, homicide, kidnapping, burglary, domestic violence and rape. Also, the records have the crime hotspots (location) and time. Furthermore, the project used one main algorithm which is the K-NN algorithm in crime identification and prediction. A prediction model was used to take the feature vector of the instance as input.

3.2 Data collection tools

3.2.1 Reports

The reports used in this project will reflect on the crimes recorded. The reports will include the crimes reported to the police by the public and also the crimes the patrol discover through other sources. the reports will be divided into two categories; the violent crimes such as homicide, domestic violence, rape and murder and property crimes such as burglary, car theft amongst other crimes involved. Using these reports will help create a the dataset because it will consist of the type of crime, where the crime has occurred as reported(location of the avenue or street) and when the crime occurred. An example from the dataset that i have homicide crime reference number that occurred in 10/31/2010 in 252 W LAKE DR NW, 105 WestLake

3.2.2 Surveys

The project will use the survey to give insights if the crime is reported or not. Officers can sometimes regard crime as no trouble because some are petty hence not considered as a major crime. Also, with the survey, it will give a clear insight to the crime statistics for example in this project we have domestic violence which can be insignificant since there are relationships involved making it hard for such a case to be reported.

3.3 System design/ Architectural framework

3.3.1 Architectural Framework

Architectural framework also known as the system design will have the user interface which connects with the user system of its Graphical User Interface nature, whereby the administrator will log in and have their details outlined. In this case since it will involve the authority, it will consist of the user administrator and the password and this will be saved in the system database. The database system will contain the log in details: administrators name and password, the k-means cluster algorithm which will highlight the state, district, type of crime, latitude and longitude. The system will require to upload the CSV

file which entails the dataset to train it using machine learning algorithm. Once the data has been uploaded, the training algorithm model will use the dataset which was imported and train it which should give the result of the prediction. Prediction in this case will highlight if its a high or low crime rate area, the number of cases involved for example the number of homicide cases in is 13, and show the analysis which can be in form of a pie chart.

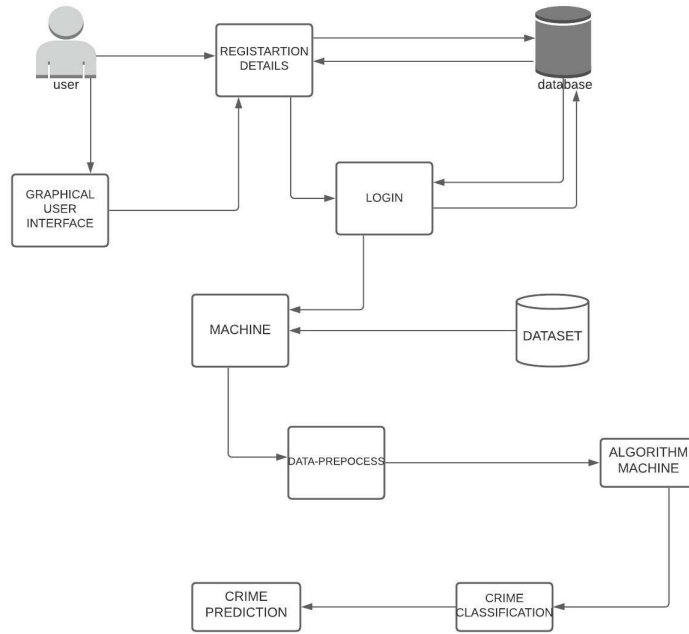


Figure 5: Architectural framework

The framework embraces the crimes committed and different types. We primarily use machine learning algorithms for prediction and find the patterns of crime. The machine learning algorithms include linear regression, k-nn means cluster algorithms as well as he apriori algorithms, but in this case we use the k-means cluster algorithm for the prediction. This technique tends to identify the crime patterns by use of the dataset that will be trained and give prediction of the crime areas, if its a high crime rate area or a low crime rate area. The modules used in this framework are as explained below:

3.4 Implementation

3.4.1 Administrator

This will be created using HTML ,CSS, JavaScript and php to perform the following functions; The administrator should monitor the login details by the user. Also, the administrator will monitor the current system in real time and

monitor the crime trends to ensure consistency . Furthermore the administrator will update the systems and access the system made

3.4.2 User Interface

Implementation of the user interface used CSS, HTML and JavaScript for the programming languages. Sublime text being used as common editor which is available in windows, linux and Mac os will be used. Furthermore, its freely available. The user interface will consist of the home page, login page, prediction and analysis.

3.4.3 Registration module

The user data is saved on Mysql database to allow the user login to the project. Once the registration is completed, the user can access the project.

3.4.4 Dataset

Implementing the dataset has been set using the records and surveys and generated by the CSV which is a file extension used in spreadsheets n this case Microsoft Excel. The CSV file as an output will be used as the dataset.

3.4.5 Prediction algorithm

The prediction algorithm implemented in this case is the K-means algorithm implemented by the Eucladian distance as a metric of similarity. To use this we will use python to initialise. We use the library Scikit-learn as the machine library, NumPy, Pandas and Scipy. .

3.4.6 Crime Database

The system database is implemented using SQL(Structured Query Language). SQL will use the Microsoft Server IDE, which helps the database management system that will help in storing, retrieving and updating data in the project.

3.5 Relationship between conceptual framework and architectural framework

The architectural model is key to ensure that the user or administrator has a good experience in terms of the navigation and this should not cross the experienced computer users , and the system should also perform executions faster and give up-date information thus making it interactive with the user. In this case, the architecture of the proposed system shows the connection in each of the modules from the user to the crime prediction creating a visualisation of the design of the proposed system. The conceptual framework on the other hand simply shows the input of the data, how its classified and the output used for prediction.

3.6 Tools

3.6.1 Python

Using python programming language will help implement a simple algorithm and predict the type of crime happening in a particular area. Training the models and algorithms will require some of the python libraries such as Numpy, Flask and Scikit-learn and for prediction as well.

3.6.2 Flask framework

Flask is one of python's libraries that will be used for the web application framework and is efficient to use compared to Django.

3.6.3 Scikit library

Scikit-learn a python library imported that will be used in making the machine learning model. It will build a model that will find a pattern that maps input data, finalize the model.

3.6.4 Numpy

This is a python library which handles multidimensional data and perform scientific and mathematical operations.

3.6.5 Pandas

This is an open source library which will provide the tools for data analysis using python . In this project, it will be used in specific machine learning algorithms.

3.6.6 Mysql

This is an open source database management system that works in many languages. It will be used as a backend for data handling whereby users can insert and retrieve data by executing queries written in SQL.

3.7 Testing

3.7.1 Administrator testing

The functionality of the module in the administrator will be tested by the developer who acts as the main administrator. Also, the administrator will test the dataset and check whether the information is accurate; in terms of the specificity of the location given (longitudes and latitude) and the type of crime committed.

3.7.2 Prediction algorithm testing

The functionality of each module in the prediction end will allow the user to run the code using the graphical user interface part whereby the user can run the code for successful prediction. The prediction model and algorithm will be used to train the data and predict the crime.

3.7.3 Crime database testing

On the database back-end testing technique, map testing will be done to ensure the database components store the right information and match the insertion fields on the user's end Functional testing will be conducted to ensure the queries run properly and the functions run from the front-end maps if the functions require the database queries.

References

- [1] H. U. Ay, A. Aysu Öner, and N. Yıldırım, “Can tech4good prevent domestic violence and femicides? an intelligent system design proposal for restraining order implementations,” in *International Conference on Electronic Government and the Information Systems Perspective*. Springer, 2021, pp. 34–45.
- [2] M. Vaquero Barnadas, “Machine learning applied to crime prediction,” B.S. thesis, Universitat Politècnica de Catalunya, 2016.
- [3] M. McGuire, “Technology crime and technology control: Contexts and history,” in *The Routledge handbook of technology, crime and justice*. Routledge, 2017, pp. 35–60.
- [4] C. Hollings, U. Martin, and A. C. Rice, *Ada Lovelace: The making of a computer Scientist*. Wiley Online Library, 2018.
- [5] S. Muggleton, “Alan turing and the development of artificial intelligence,” *AI communications*, vol. 27, no. 1, pp. 3–10, 2014.
- [6] F. Lara, “Artificial neural networks: An introduction,” *Instrumentation and Development*, vol. 3, no. 9, 1998.
- [7] L. E. Peterson, “K-nearest neighbor,” *Scholarpedia*, vol. 4, no. 2, p. 1883, 2009.
- [8] G. C. Oatley and B. W. Ewart, “Crimes analysis software: ‘pins in maps’, clustering and bayes net prediction,” *Expert Systems with Applications*, vol. 25, no. 4, pp. 569–588, 2003.
- [9] T. Cheng and M. Adepeju, “Detecting emerging space-time crime patterns by prospective stss,” in *Proceedings of the 12th international conference on geocomputation*, 2013.
- [10] D. J. Fitzpatrick, W. L. Gorr, and D. B. Neill, “Policing chronic and temporary hot spots of violent crime: A controlled field experiment,” *arXiv preprint arXiv:2011.06019*, 2020.
- [11] N. Shah, N. Bhagat, and M. Shah, “Crime forecasting: a machine learning and computer vision approach to crime prediction and prevention,” *Visual Computing for Industry, Biomedicine, and Art*, vol. 4, no. 1, pp. 1–14, 2021.
- [12] R. Hartman, “Ai is watching you.”
- [13] H. Yu, L. Liu, B. Yang, and M. Lan, “Crime prediction with historical crime and movement data of potential offenders using a spatio-temporal cokriging method,” *ISPRS International Journal of Geo-Information*, vol. 9, no. 12, p. 732, 2020.

- [14] X. Zhang, L. Liu, L. Xiao, and J. Ji, "Comparison of machine learning algorithms for predicting crime hotspots," *IEEE Access*, vol. 8, pp. 181 302–181 310, 2020.
- [15] N. Abdulrahman and W. Abedalkhader, "Knn classifier and naïve bayse classifier for crime prediction in san francisco context," *International Journal of Database Management Systems (IJDMS)*, vol. 9, no. 4, pp. 1–9, 2017.
- [16] A. Almaw and K. Kadam, "Survey paper on crime prediction using ensemble approach," *International Journal of Pure and Applied Mathematics*, vol. 118, no. 8, pp. 133–139, 2018.
- [17] J. Agarwal, R. Nagpal, and R. Sehgal, "Crime analysis using k-means clustering," *International Journal of Computer Applications*, vol. 83, no. 4, 2013.
- [18] R. Kiani, S. Mahdavi, and A. Keshavarzi, "Analysis and prediction of crimes by clustering and classification," *International Journal of Advanced Research in Artificial Intelligence*, vol. 4, no. 8, pp. 11–17, 2015.
- [19] X. Chen, Y. Cho, and S. Y. Jang, "Crime prediction using twitter sentiment and weather," in *2015 Systems and Information Engineering Design Symposium*. Citeseer, 2015, pp. 63–68.