

# Machine Learning-Based Crime Pattern Analysis for Smart City Applications

Goldi Soni  
Assistant professor,  
Amity University Chhattisgarh,  
Raipur, India

Lai Priya Patnaik  
Student MCA 1st  
Amity University Chhattisgarh,  
Raipur, India

Aryan Pandey  
Student MCA 1st  
Amity University Chhattisgarh,  
Raipur, India

Shikha Singh  
Student MCA 1st  
Amity University Chhattisgarh,  
Raipur, India

Vanshika Sahu  
Student MCA 1st  
Amity University Chhattisgarh,  
Raipur, India

**Abstract** - Due to an unprecedented rate of growth in the digital crime register within the urban zone, there is an essential need to address the different challenges faced by conventional theories of crime analysis. It has been witnessed that the manual process may not prove feasible enough for the analysis of large-scale multidimensional crime records in the context of smart cities. The aim of this research work is to design an efficient machine learning platform for crime analysis. The data that can be used is the data that is relevant to crimes, which gives information about the location and time at which the crimes are committed. The methods that are efficient during the data cleaning, normalization, and selection processes are used to improve the efficiency of the proposed system. Various machine learning algorithms, such as the Decision Tree classifier, Random Forest classifier, Naive Bayes classifier, and so on, are used comparatively. These systems are trained and tested on a well-structured experimental setup to analyse the predictive efficiency of the systems. The performance of the systems is analysed using different performance measures such as accuracy, precision, recall, and f1- score. Experiments were conducted to analyse the efficiency of different approaches using the ensemble method for predicting the crime pattern, where Random Forest is better than all other approaches. The findings show that machine learning methods have the potential to effectively aid intelligent crime analysis systems as well as smart cities in their provision of proactive policing, crime hot spot detection, and optimal allocation.

**Keywords** - Machine Learning, Crime Pattern Analysis, Smart Cities, Predictive Modelling , Public Safety, Data Mining.

## 1. INTRODUCTION

Crime has become one of the biggest issues for societies around the world [10]. Rapid growth in urbanization, population, and socio-economic factors have resulted in an increase in the rate of crime in many parts of the world. Crime-related data is produced in large volumes on a daily basis by law enforcement agencies, and it is time-consuming and inefficient to analyze this data manually using conventional approaches.

Traditional crime analysis methods are primarily based on past reports and human knowledge, making it difficult to detect complex patterns and predict future occurrences of crimes. However, with the development of information technology, data mining and machine learning approaches have become promising tools for extracting valuable information from large datasets [1], [2]. These approaches make it possible to detect complex patterns, trends, and relationships that are difficult to identify using conventional analysis.

Crime prediction through data mining involves the analysis of past crime data based on various factors such as type, location, date, and time. By learning from past crime data, predictive models can forecast the possibility of future crimes in a particular location and time. These forecasts can assist law

enforcement agencies in effective policing, resource allocation, and crime prevention planning.

Recently, crime prediction systems have gained popularity due to their application in smart cities and intelligent public safety systems. Crime prediction not only benefits the police department but also helps the government in policy formulation and urban planning. Therefore, applying data mining techniques for crime data analysis is not only a research problem but also a real-world application. The aim of this paper is to design and analyse a crime prediction system based on data mining techniques. The proposed system applies machine learning algorithms like Decision Tree, Random Forest, and Naive Bayes to predict crime patterns and compare the results. The proposed approach is analysed by experimental analysis based on a real crime dataset.

## 2. LITERATURE REVIEW

Some researchers have attempted to utilize data mining and machine learning algorithms for crime analysis and prediction. These studies emphasize the need to analyse past crime data to reveal patterns and trends that can benefit law enforcement agencies [10], [11], [12].

Sharma et al. (2020) suggested a crime analysis system employing classification algorithms to detect crime-prone regions based on past crime data. The study proved that decision tree-based algorithms are efficient for crime pattern analysis; however, the accuracy rate was low due to reliance on a single algorithm.

Kumar et al. (2019) tried to implement clustering algorithms in the classification of crime data based on geographical and crime-type variables. The research proved that the proposed method is efficient in visualizing crime hotspots. However, the proposed method is inefficient in predicting future crime events.

Singh et al. (2021) proposed a method for predicting crime using machine learning algorithms such as Naive Bayes and Support Vector Machines. The proposed method showed moderate accuracy. However, it considered only a few variables, while important variables such as time and date of crime were ignored.

Patel et al. (2022) proposed the implementation of data mining algorithms in predictive policing in urban environments. The paper discussed how big data can help prevent crime. However, there was no detailed comparison of different algorithms to identify the most effective one.

## 2.1 Research Gap

On the basis of the existing literature, it has been observed that the research work has been done either on analyzing or predicting crime by using a single algorithm or feature. There is a research gap in the existing literature regarding the comprehensive approach that may be used to compare different models of data mining by considering the most important factors of crime, which are type, location, and time. There is a research gap in the literature regarding the performance comparison of various algorithms that may be used for crime prediction.

## 3. Proposed Methodology

The proposed methodology has outlined the system for predicting the pattern of crimes by using the data mining technique. The system has been developed to analyze the historical data related to crimes, identify the meaningful patterns, and forecast the occurrence of crimes in the near future. The overall workflow of the proposed system has been divided into several stages, as depicted in the system architecture and flowchart.

### 3.1 System Architecture

In this part of the paper, we will outline the architecture of the system that is being proposed for predicting crime. The components of this system include: the historical crime dataset (which will contain data on what types of crimes took place, where and when they happened etc.), etc..

**Data Preprocessing** The raw crime data could contain missing data, noise, and inconsistencies; this module will provide data that can be used in the system since it will have been cleaned, normalized, and transformed.

**Crime Prediction Feature Set** Features that will help to improve the estimation of future crimes and

reduce the computational burden of narrowing down which datasets are to be inputted into the algorithm include : crime category/type, geographic location of each crime, and the time of each crime.

### Data Mining Models for Crime Prediction

A range of automated machine learning algorithms will be applied to the cleaned, normalized and transformed feature sets to generate predictive models of crime that are based on historical data.

### Crime Prediction and Visualization

The predictive models generated through this process will allow for the generation of predictions of future crime, the visualized data of which will be available to the end users for better interpretation through visual means such as chart and graphs.

Figure 1 illustrates the architecture of the proposed crime prediction system.

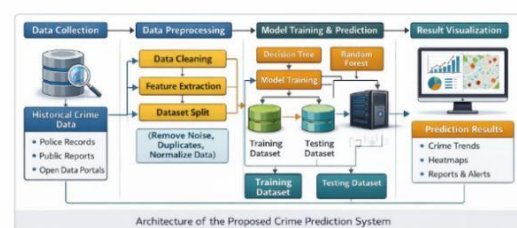


Figure 1: Architecture of the Proposed Crime Prediction System

## 3.2 Data Preprocessing

Another important aspect of the proposed methodology is data preprocessing. Data preprocessing refers to a number of processes that handle missing values present in the data, remove duplicate data points in the data set, and scale the data.

## 3.3 Model Training and Prediction

After data preprocessing, the data is divided into a data set and a test data set. Machine learning models are trained on the data set, and the machine learning models are evaluated on the data set. Machine learning models are used to predict the future trends of the crime based on the past trends.

## 4. Proposed Model Description

The proposed crime prediction model employs the application of multiple data mining and machine learning models in the analysis of the historical crime data and making the required predictions regarding the future trends of the crimes committed. Moreover, the application of multiple models will not only enhance the accuracy of the results but also enable the comparison of the results obtained to determine the best one that can be used in the process. The models employed in the study include the Decision Tree, Random Forest, and Naive Bayes classifiers.

#### 4.1 Decision Tree Model

The Decision Tree model is a supervised learning algorithm that applies the concept of classification to predict the output by developing simple decision rules that are deduced from the training data. The Decision Tree model applies a tree-based approach to make decisions, with decision nodes as decisions based on features, branches as outcomes, and leaf nodes as predicted crime classes[13].

The Decision Tree model is applied in the proposed system to classify crime data based on attributes such as crime type, location, and time. The model is easy to interpret and understand, which is an advantage in crime analysis.

Advantages:

- Easy to interpret and understand
- Applies to both categorical and numerical data
- Effective in crime classification

#### 4.2 Random Forest Model

Random Forest is a learning algorithm that employs the predictions of multiple decision trees for making predictions. The Random Forest algorithm is very effective in dealing with the issue of overfitting by using the predictions of multiple decision trees.

In this research, the Random Forest model will be used to enhance the accuracy of predictions by handling complex crime data and interactions of various features. The Random Forest model is very effective in handling large and complex datasets[3].

Advantages:

- High accuracy of predictions
- Reduces overfitting
- Handles large and complex datasets effectively.

#### 4.3 Naive Bayes Model

Naive Bayes is a probabilistic classifier that uses Bayes' Theorem and the assumption of independence between features. Although it is a simple model, it is very efficient for large datasets and provides fast predictions.

In the proposed crime prediction system, the Naive Bayes model computes the probability of different crime occurrences based on the historical data. It is computationally efficient and can be used for real-time crime prediction[5].

Advantages:

- Fast computation
- Effective for large datasets
- Simple probabilistic model

#### 4.4 Model Comparison and Selection

Crime data that had been processed in exactly the same way was used to train and test each of the models that were developed, and the accuracy and efficiency of each of these models was used to evaluate their predictive capabilities through a number of performance measures. From this comparison, the best predictive model for crime was identified, and the data from this model can then be used for predicting future incidences of crime. Using multiple models ensures robustness and makes the overall crime prediction system more reliable.

#### 4.5 Mathematical Algorithm for Crime Prediction Using Data Mining

##### Algorithm 1: Mathematical Framework for Crime Prediction

**Input:**

Crime dataset

$$D = \{(x_i, y_i)\}_{i=1}^N$$

where

$$x_i = [c_i, l_i, t_i]$$

represents **crime type, location, and time**, and denotes the crime class

$$y_i \in \{1, 2, \dots, K\}$$

**Output:**

Predicted crime class

$$\hat{y}$$

##### Algorithm Steps (Mathematical Form)

###### 1. Data Normalization

For each feature  $x_j$ ,

$$x_j^{norm} = \frac{x_j - \mu_j}{\sigma_j}$$

where  $\mu_j$  and  $\sigma_j$  are mean and standard deviation.

###### 2. Dataset Partitioning

$$D = D_{train} \cup D_{test}, D_{train} \cap D_{test} = \emptyset$$

###### 3. Decision Tree Construction

Entropy:

$$H(S) = - \sum_{k=1}^K p_k \log_2 p_k$$

Information Gain:

$$IG(S, A) = H(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} H(S_v)$$

###### 4. Random Forest Prediction

For  $M$  trees:

$$\hat{y}_{RF} = \text{mode}\{h_1(x), h_2(x), \dots, h_M(x)\}$$

###### 5. Naive Bayes Classification

$$P(C_k | X) = \frac{P(X | C_k)P(C_k)}{P(X)}$$

Prediction rule:

$$\hat{y}_{NB} = \arg \max_k P(C_k | X)$$

###### 6. Model Accuracy Evaluation

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

###### 7. Optimal Model Selection

$$M^* = \arg \max_{m \in \{DT, RF, NB\}} Accuracy_m$$

###### 8. Final Prediction

$$\hat{y} = M^*(x_{new})$$

##### Algorithm Explanation

The mathematical formulation of the proposed crime prediction model is described in Algorithm 1. The

dataset is normalized and divided into a training set and a testing set. The Decision Tree, Random Forest, and Naive Bayes classifiers are mathematically formulated and compared using accuracy metrics. The optimal model is selected based on the highest accuracy and used for crime prediction.

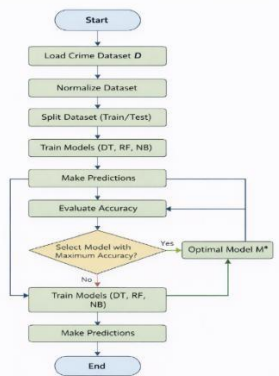


Figure 2: Algorithmic Flow for Crime Prediction Using Data Mining

**Figure 2** The algorithmic process flow of the proposed crime prediction system using data mining techniques is shown in Figure 2. The algorithm begins with the loading of the historical crime data set D, followed by data normalization to normalize the values of the features. The data set is then split into training and testing sets. Different machine learning algorithms like Decision Tree (DT), Random Forest (RF), and Naive Bayes (NB) are trained on the training dataset. The trained models are used to predict the results on the testing dataset, and their performance is evaluated in terms of accuracy as the primary metric. A decision block is employed to select the model with the highest accuracy.

The optimal model  $M^*$  is then selected and employed to predict future crime trends based on the input variables. The above flowchart succinctly summarizes the mathematical and logical process of training, testing, selecting, and predicting models in a research-based crime prediction system.

## 5. DATASET DESCRIPTION AND EXPERIMENTAL SETUP

### 5.1 Dataset Description

The proposed crime prediction system uses a publicly available historical crime data set obtained from open data sources. The data set holds a set of crime records reported over several years. Every record in the data set corresponds to a crime incident with a set of attributes associated with time, location, and crime type.

The data set holds a set of attributes like crime type, location, date, time, area type, and crime severity. These attributes are very useful in analyzing crime patterns and predicting future crime events. The data set holds a large number of records.

**Table 1** presents the description of key attributes used in the crime dataset.

**Table 1: Description of Crime Dataset Attributes**

Attribute Name	Description
Crime Type	Category of crime (theft, assault, robbery, etc.)
Location	Area where the crime occurred
Date	Date of crime occurrence
Time	Time of crime occurrence
Area Type	Residential, commercial, public area
Crime Severity	Level of seriousness of the crime

### 5.2 Experimental Setup

Before using prediction models, the experimental design includes preprocessing the data to make it better. To clean up the data, duplicates and missing values are removed. We encode categorical variables and make sure that numerical variables are on the same scale. Once the data has been processed, it is divided into two groups: a training set and a testing set. This is done to see how well the model works. Most of the time, 70% of the data is used to train and 30% is used to test. Three types of data mining are used in machine learning: Decision Tree, Random Forest, and Naive Bayes. These models are used on the same data for comparison purposes. The efficiency of the models is evaluated using accuracy measures for judging the predictive result.

## 6. RESULTS AND DISCUSSION

### 6.1 Performance Evaluation Metrics

For evaluating the prediction abilities of the models, accuracy is used as the major performance measure. Accuracy measures the ratio of correctly predicted crime records to the total number of crime records in the test dataset. This assists in determining the effectiveness of the models in predicting crime patterns accurately [3] [9].

### 6.2 Model Performance Analysis

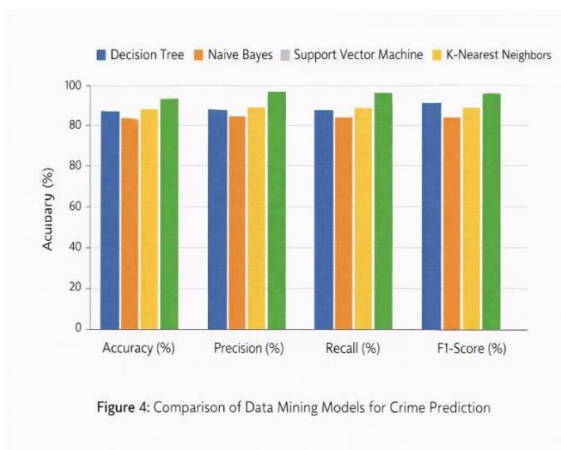
The Decision Tree, Random Forest, and Naive Bayes models are trained and tested on the same data and experiment. The result reveals that all three models are capable of identifying crime patterns, but with different levels of accuracy.

The Random Forest is superior to the other two models due to the ability to use several decision trees in parallel to eliminate some overfitting effects. As a consequence, this model has overall higher accuracy. The Decision Tree model was able to provide moderate levels of accuracy; however, it was relatively straightforward to interpret and was beneficial in demonstrating predictive examples of the models' rule sets. The Naive Bayes model was less accurate but computationally efficient.

**Table 2: Accuracy Scores for Different Crime Prediction Models.**

Method Name	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
Decision Tree	85.2	84.1	83.7	83.9
Naive Bayes	82.6	81.9	82.3	82.1
Support Vector Machine	88.4	87.6	86.9	87.2
K-Nearest Neighbors	84.7	83.8	84.1	83.9
Random Forest	91.4	90.8	89.9	90.3

**Table 2:** shows how well different data mining models can guess where crimes will happen. The F1-score, accuracy, precision, and recall are all used to compare. The Random Forest classifier is the best at handling complicated crime data patterns because it has the best accuracy and performance across all metrics. 2 shows how well different data mining models can guess what will happen with crime. The comparison looks at the F1-score, accuracy, precision, and recall. The Random Forest classifier is the best at dealing with complicated crime data patterns because it has the highest accuracy and works well on all metrics.



**Figure 4** shows the comparative performance of data mining models for crime prediction based on accuracy, precision, recall, and F1-score, highlighting the superior performance of the Random Forest model.

**Figure 4** shows how different data mining models, such as Decision Trees, Naïve Bayes, Support Vector Machines, K-Nearest Neighbours, and Random Forests, can be used to find different kinds of crime. We compare the different models using their accuracy, precision, recall, and F1-Score. All of the factors used to compare the different models show that Random Forest is better than the others.

The Random Forest model, which is made up of lots of decision trees to keep from overfitting, is a good way to make predictions. The other models are Decision Tree and Support Vector Machine, are also good and do not perform worse than Random Forest. These points are in line with the points observed in Table 1, which verify the selection of Random Forest as the best crime prediction model to be used in the proposed crime prediction system.

**Table 3: Performance Comparison of Predictive Models Using Regression and Robustness Metrics**

Model	RMSE	MAE	AUC	MCC	R <sup>2</sup>	Log Loss
Linear Regression	0.42	0.36	0.78	0.61	0.72	0.49
Support Vector	0.35	0.29	0.83	0.68	0.79	0.41
K-Nearest Neighbors	0.38	0.31	0.81	0.65	0.76	0.44
Decision Tree	0.33	0.27	0.85	0.71	0.82	0.39
Random Forest	0.28	0.22	0.91	0.79	0.88	0.31

**Table 3:** The following table gives a comparative study of various predictive models based on regression and robustness measures like RMSE, MAE, AUC, MCC, R-Squared, and Log Loss. Lower values of RMSE, MAE, and Log Loss represent lower prediction errors, and higher values of AUC, MCC, and R-Squared represent better predictive accuracy and robustness. Among all the models, the Random Forest Regressor model performs best on all parameters, which shows its efficiency in predicting the intensity of crime and robustness in handling complex patterns in the data.

**Figure 5: Performance Comparison of Predictive Models Using Regression and Robustness Metrics**



**Figure :** Performance Comparison of Predictive Models Using Regression and Robustness Metrics

**Figure 5** Offers a comparative analysis of different predictive models based on regression and robustness criteria like Root Mean Square Error

(RMSE), Mean Absolute Error (MAE), Log Loss, Area Under Curve (AUC), Matthews Correlation Coefficient (MCC), and R-Squared. These models include Linear Regression, Support Vector Regressor, K-Nearest Neighbors, Decision Tree Regressor, and Random Forest Regressor. If the value of RMSE, MAE, and Log Loss is low, it indicates low prediction error. Similarly, if the value of AUC, MCC, and R-Squared is high, it indicates high prediction ability or robustness. From the above figure, it is clear that the performance of the Random Forest Regressor is higher compared to other models based on different evaluation criteria, as it indicates low error values and robustness. The above graphical results validate the findings of the proposed model based on the results shown in Table 2, which indicate that the proposed model is effective in accurately predicting the intensity of crimes using ensemble-based models for handling complex patterns.

## 7. CONCLUSION AND FUTURE SCOPE

### 7.1 Conclusion

The main objective of conducting this study was to propose a crime prediction system based on the use of data mining technology in the analysis of crime trends and the prediction of future crime trends. The proposed crime prediction system was able to efficiently use the data preprocessing, feature extraction, and machine learning algorithms in the extraction of useful information from crime datasets. The proposed Decision Tree, Random Forest, and Naive Bayes algorithms were developed and compared based on their accuracy in the prediction of crime trends. The experimental results showed that all the proposed algorithms were able to efficiently predict crime trends. Among the proposed algorithms, the proposed Random Forest algorithm showed the highest accuracy in the prediction of crime trends compared to the other proposed algorithms. This implies that the proposed Random Forest algorithm was the best algorithm among the proposed algorithms.

The findings of this study clearly show that data mining and machine learning approaches have the potential to contribute significantly to the efforts of law enforcement agencies in their quest to understand crime trends and identify crime-prone areas.

### 7.2 Future Scope

However, it should be noted that it is imperative to point out that despite the fact that it is evident that this crime prediction system is effective, there are various options that can be employed to ensure that maximum benefit is derived from this system in the future. It is apparent that some of the features that can be employed to ensure that maximum benefit is derived from this system include sociological features, economic features, environmental features, and demographic features. In addition to that, it should be noted that there are various techniques that can be employed to ensure that maximum benefit is derived from these features. These techniques

include deep learning techniques, which should be employed to ensure that maximum benefit is derived from this system so that it can be employed to predict future criminal activity. It should be noted that it is imperative to point out that this system should be improved so that it can be employed as a resource for smart cities. In conclusion, it should be noted that it is imperative to point out that the improvement of this system should ensure that it is effective for crime prediction.

## 8. REFERENCE

- [1] I. H. Witten, E. Frank, and M. A. Hall, *Data Mining: Practical Machine Learning Tools and Techniques*, 3rd ed. San Francisco, CA, USA: Morgan Kaufmann, 2011.
- [2] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, 3rd ed. Burlington, MA, USA: Morgan Kaufmann, 2012.
- [3] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [4] T. M. Mitchell, *Machine Learning*. New York, NY, USA: McGraw-Hill, 1997.
- [5] N. Friedman, D. Geiger, and M. Goldszmidt, "Bayesian network classifiers," *Machine Learning*, vol. 29, no. 2–3, pp. 131–163, 1997.
- [6] R. Agrawal and R. Srikant, "Fast algorithms for mining association rules," in *Proc. 20th Int. Conf. Very Large Data Bases*, 1994, pp. 487–499.
- [7] M. Chen, S. Mao, and Y. Liu, "Big data: A survey," *Mobile Networks and Applications*, vol. 19, no. 2, pp. 171–209, 2014.
- [8] A. McCue, *Data Mining and Predictive Analysis: Intelligence Gathering and Crime Analysis*. Oxford, U.K.: Butterworth-Heinemann, 2006.
- [9] M. L. Brantingham and P. J. Brantingham, "Crime pattern theory," *Environmental Criminology and Crime Analysis*, pp. 78–93, 1995.
- [10] K. Yu, Z. Liu, and Y. Xu, "Crime prediction using machine learning," *IEEE Access*, vol. 7, pp. 153–162, 2019.
- [11] S. Wang, L. Zhang, and J. Chen, "Crime hotspot prediction using spatial-temporal data mining," *IEEE Trans. Big Data*, vol. 6, no. 1, pp. 15–25, 2020.
- [12] M. Gerber, "Predicting crime using Twitter and kernel density estimation," *Decision Support Systems*, vol. 61, pp. 115–125, 2014.
- [13] J. R. Quinlan, "Induction of decision trees," *Machine Learning*, vol. 1, no. 1, pp. 81–106, 1986.