# AN ANALYSIS OF CRIME PREDICTION USING MACHINE LEARNING AND DEEP NEURAL NETWORKS

[1]Ajeet Singh [2] Kunal Sharma [3]Md Vadiyat Naqvi [4]Khushi Tyagi [5]Kumkum Joshi

[1] Asst. Prof, Dept of CSE, Moradabad Institute of Technology, Moradabad UP, India
ajeetsingh252@gmail.com

[2][3][4][5] Student, Dept of CSE, Moradabad Institute of Technology, Moradabad UP, India
kunalsharma8630@gmail.com; mdvadiyatnaqvi@gmail.com;
khushityagi771@gmail.com; kumkumjoshi066@gmail.com

***ABSTRACT***

*Loss of life and property damage can be prevented by preventing crimes. Using machine learning to prevent crime is received a lot of research. This article reviews the most advanced crime prediction techniques that have become available in recent years 10 years, discuss possible challenges and propose possible future trends. Many articles have tried to predict crimes, but datasets and methods used in different applications.*

*Using the Systematic Literature Review (SLR) technique, we aim to collect and synthesize the necessary knowledge on machine learning based crime prediction to help both law enforcement agencies and investigators reduce and prevent future crimes. Our main focus is on 68 selected machine learning studies that predict crime. We select eight research areas and find that despite the fact that labeled data may not always exist in real situations, most papers used supervised machine learning techniques assuming that labeled data existed in the past.*

*We have also looked at the main issues that the researchers faced during the course of some of these studies. We believe that this research opens up new avenues of research to help countries reduce crime rates for improved safety and security.*

***KEYWORDS:*** *Machine Learning, Crime Prediction, Supervised Machine Learning, Neural Network*

## 1. INTRODUCTION

A crime is a physical attack or an illegal act perpetrated against another person by an individual or a group of individuals against another person. A crime may result in physical injury or damage to the property of the victim. A crime is prohibited by the laws of the place where it occurred. Law enforcement agencies use crime solving tools but often fail to implement effective prevention measures. Crime has been on the rise in all countries.

In another study, Mattereke (2021) ranked South Africa, Venezuela, and Papua New Guinea (PNG) as the top 3 countries most prone to crime. Walczak (2021) notes that as a country's population grows, so does the crime rate. Crime prediction becomes less reliable as a result of population growth (Pratibha (2020)). A country's security and safety are important for its economy and social wellbeing. Preventing crime will reduce economic losses while enhancing public safety (TopiReddy (2018) and Rumi (2018).

Data mining is the process of extracting statistical concepts, patterns, and correlations from massive amounts of data with no user experience and subjectivity. Criminal detectives can use this powerful tool to easily and effectively analyze large data sets. Machine Learning and Data Mining are both complex topics that need computers, mathematics, and programming in order to perform certain tasks. These fields are essential for preventing and detecting crime. Data Mining is the process of discovering new patterns in massive amounts of data by using statistical, artificial techniques intelligence, and database management. Clustering is mostly utilized to analyses crime hotspots in crime-related research, primarily by grouping crime sites or time series flows [2,3], while the generated patterns of clustering

are used to discover crime trends [4]. For crime prediction, various approaches have been developed. Crime impedes a nation's and its communities' growth. Crime is a widespread epidemic that exacerbates social and economic problems (Mattereke et al. 2021). Characteristics, sometimes called social and economic factors (S&E factors), determine whether a criminal commits or suffers a crime or not. In many countries, S&E factors help to reduce crime rates by predicting when crimes will occur (see Hajela et al., 2020).

Many researchers have been motivated to continue their research in the area of crime due to the large amount of data published by some governments. Several researchers have developed several models to predict the likelihood of future crimes due to historical data which makes it an interesting subject matter for research. Due to limitations on its use law enforcement organizations may sometimes refuse to provide their data to local researchers which further exacerbates the frustrations and disappointments. Police and government have a large amount of data which can be used to reduce crime rates. The crime pattern theory states that offenders prefer to commit opportunistic or violent crimes by using places they are familiar with, rather than going into unknown territory. (see Jalil et al., 2017)

Deep Learning, also known as Machine Learning or ML, is a branch of AI that is currently used to make more intelligent decisions and predict future events in many areas. ML is defined as the study of computer algorithms which automatically improve themselves by using data, experience, and learning. Deep learning (DL) is inspired by the way our brain works. It is a type of artificial neural network that uses multiple layers and layers types (such as pooling, convolution, fully connected, dropout, etc.) in an effort to mimic how our brain works. There are four kinds of learning: reinforcement learning unsupervised learning semi super supervised learning supervised learning crime prediction model. To develop a highly accurate crime prediction model, it is essential to comprehend the nature of a crime (Elluri et al. 2019). The age, gender, location, number of perpetrators, education level, income, type of weapon used, and victims age, gender, location, economic status, and educational level are just a few examples of aspects that can be included in a crime's characteristics.
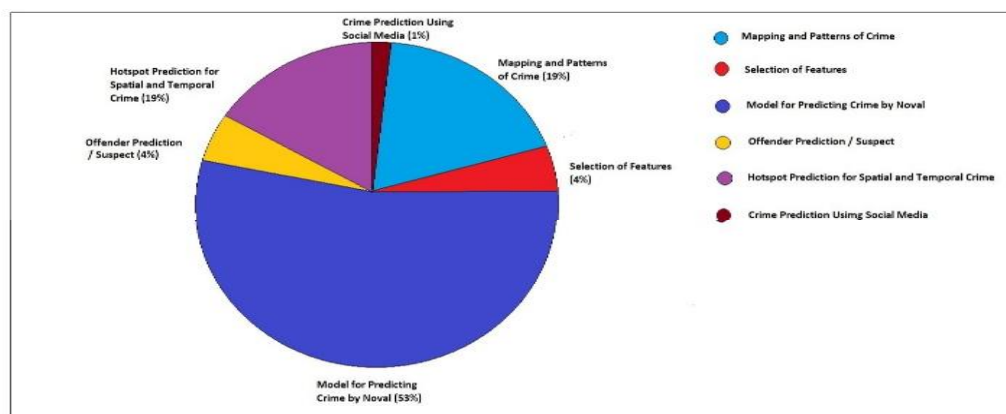


**Figure 1:** Distribution of Objectives

## 2. BACKGROUND AND HISTORY

In criminology, crime control and suppression are the two main categories of data mining. De bruin et al. developed a crime trend model that categorizes people based on how well their profiles match each other. Manish Gupta describes the technology that Indian police are currently using to achieve their e-governance objectives. They urge the police to adopt an interactive query-based crime analysis interface. He created an interface that extracts essential information from the massive crime database held by the NCRB using crime data mining techniques such as cluster and others. In the case of Indian crime records the proposed interface has proven to be very useful. Sutapat THIPRUNGSRI examines the use of cluster analysis in accounting, specifically the detection of audit anomalies.

Residential areas near schools have a higher rate of crime than other neighbourhoods, according to a study by Rougout Roman (5). This contrasts with a study by Kaut et al. (6), which found that proximity to schools had an effect on the neighbourhood's crime rate. In Engstad's (7), a comparison was made between the number of 'similar types of crime' in a neighbourhood (with or without hotels) and the number in a neighbouring neighbourhood (without hotels). Hotels are often regarded as 'criminal facilities', as they tend to attract odd or intoxicated people, who are easy targets for theft. Blue and Davis (8) found a crime hotspot about a 15-minute walk from Chicago's tube stations, and a similar pattern was found near Bronx tube stations.

## 3. AIM AND SCOPE

Data science and machine learning will be used in this project to predict crimes using crime data sets. Crime data is collected from the police department's official web portal. Crime data includes information on the crime, such as time, date, location, and type of crime. Data will be preprocessed before training the model. Feature selection and scaling will then be performed to ensure high accuracy.

Analysing crime classification using ANN and various other techniques (decision tree and random forest) and proposing a better query based training algorithm.

The dataset will be illustrated using visual representations of hundreds of incidents, for example, at peak times of day or during peak months of the year. The main objective of this project is to illustrate just one of the ways law enforcement organizations could use machine learning to detect, predict, and solve crimes much faster and thus reduce the crime rate. Depending on the availability of the dataset, this can be implemented in different states or countries.

## 4. METHODOLOGY

Artificial neural networks (ANNs) were originally designed as machine learning techniques. The first ANN was proposed by Warren McCulloch, Walter Pitts, and others in the early 1980s as a tool to perform complex calculations like propositional logic using a computer model of biological neurons. The first ANN didn't solve XOR, and ANNs weren't widely discussed until the 1990s, when support vector machines (SVM's) and other powerful machine learning algorithms began to be developed. The invention of deep neural networks (DNNs) enabled the management of large amounts of data that wasn't being analyzed. The DNN is a model of an artificial neural network that has many layers, including the input layer and the output layer. The output layer is the weight of each node in the DNN. This allowed the prediction model to derive outputs from weights.

$$n_k^h = \sum_{j=1}^{R} w_{k_j}^h p_j + b_h^k k = 1 \ to \ S$$

**Figure 2** illustrates the ANN's basic composition,
While Figure 2 shows its equation. "P" stands for "put variable," "R" for "number of input variables," "S" for "hidden neurons," "b" for "hidden layer," and "w" for weight. The activation function takes in the weight of each
element as its input. The result is produced by adding these weighted values together. In early investigations, the sigmoid function was commonly used as an activation function. The function exhibits a phenomenon known as the gradient vanishing, wherein the present data converges to zero as the network grows.. Additionally, extra computing time is required because the sigmoid function is an exponential function. The issue was solved with the use of a nonlinear ReLU function, which does not lose information when the input value exceeds the threshold value. This method calculates the value of the gradient with a simple value of 0 or 1. Since it significantly improves ANN performance, the ReLU function has been used in various research and is also used in the current work [12]. The structure of

the hidden layers and nodes for producing the best ANN model does not have predefined standards or rules.

It's important to carry out multiple trials in order to find the best model with the MSE value. A model should have at least one hidden layer and a maximum of seven[13]. To identify the best model with the lowest MSE value, the minimum, median, and maximum hidden node counts in three to five cases had to be examined. In order to compare the model performances, the environmental factors—which serve as input data—are constructed differently in this study.
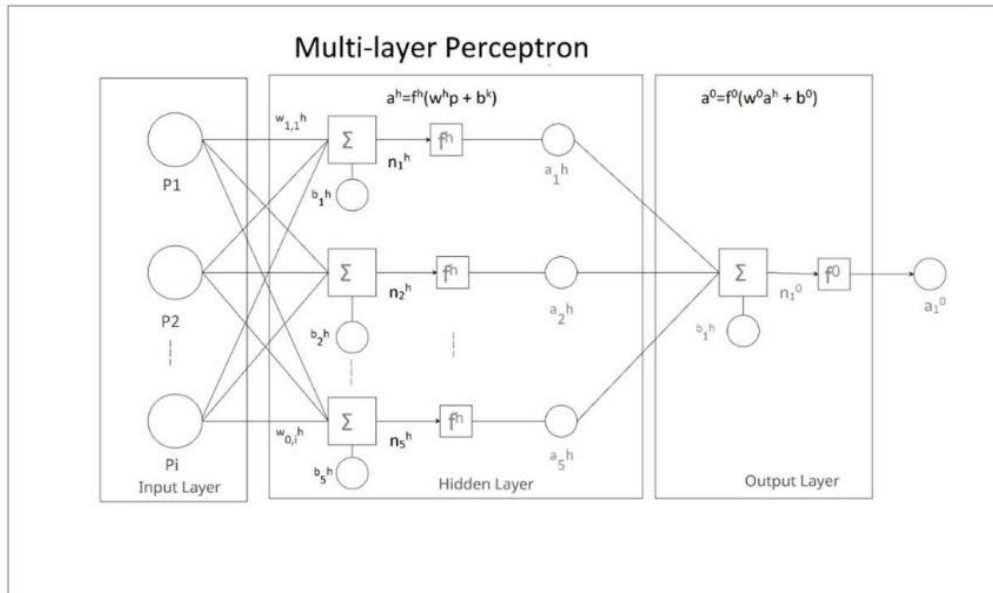


**Figure 3:** A Multi-Layer Perceptron

## 5. RESULTS AND DISCUSSION

### 5.1. Crime Categorization

In the Dongjak District, K-modes clustering indicates that there are typically four types of theft. The statistics show the different categories of theft. Cluster A is made up of other theft (53%) in business settings (43%), which primarily targets men (76%) at night (62%). Intrusion (16%), striking (19%), car theft (8%), trick (4%), and other (53%) are the most common methods of theft. Additionally, hitting theft and other stealing ratios were higher than those in other clusters, and evening crime rates (62%) were three times higher than those at other periods.

According to Hwang [14], burglars frequently choose female victims on purpose. The most frequently cited reason for choosing to target women was due to ease of control. This suggests that they believe that women are easy targets even if they have been in contact with their victims. The findings show that gender has a significant influence on the choice of crime targets.5.2.    Derivation of Significant Environmental Factors

Cohen's f2 is a common statistic for assessing various forms of impact magnitude when using F-test for ANOVA and multiple regression analysis. Cohen's f2 is frequently used in social science statistical analysis. R2 = 0.02 should be viewed as a small effect size, 0.13 as a medium effect size, and 0.26 as a large effect size, according to Cohen [15]. All of the values for the multiple regression analysis models in this study fell within the range of 0.16 and 0.28, which shows that the models have undergone appropriate validation.

Additionally, there were numerous structures in Dongjak District's commercial sections that combined residential and commercial functions. According to Kim et al. [16], this combination increased the

likelihood of crime. Because medical facilities are places where strangers congregate, Stucky and Ottensmann [17] reached to a conclusion that increased chances of severe criminal offences there.

## 6. CONCLUSION

Considering the contextual characteristics at the scene of the crime, an artificial neural network model was developed to predict the sorts of theft that are most likely to happen in random locations. The theft offenses that occurred in the Dongjak District between 2004 and 2015 were looked at in this research.. The criminal theft in the Dongjak District was first divided into four categories using k-modes clustering. Additionally, it was shown through the use of multiple linear regression analysis that each classification had distinct associated environmental

components. Based on the types of theft, it was shown that the effects and significance of environmental factors differed. Based on the forms of theft, the same environmental elements had varying degrees of impact; these findings can be used to determine the best tactical actions to avoid each type of theft. Using the prediction model developed during the study's crime prevention stage, we could predict the types of theft offenses that are most likely to occur in a specific place and identify the environmental factors that affect the crime type's likelihood. This allows for the creation of an effective crime prevention strategy by identifying an environmental component that has to be modified primarily for each type of crime. This allows for the creation of an effective crime prevention strategy by identifying an environmental component that has to be modified primarily for each type of crime.

When utilizing the hybrid algorithm to predict the crime that occurs in a certain location in a single dataset, this research achieved a decent result. In order to get the best results on a huge dataset, ANN gives the system the capacity to learn from both past and present events utilizing SOM and K-means. With the help of a digital representation provided by GIS, users are able to map crime scenes analytically and descriptively from their various geographic locations and to visualize data in a variety of ways that highlight relationships. This data can also be divided into layers using GIS technology. In this thesis, a methodology for creating a hybrid model approach has been described. The approaches are applied to actual data on the distribution of burglary incidence in the studied area. For the analysis of the observed data, relevant statistical concepts, GIS, clustering algorithms, and ANN are applied.

## REFERENCES

[1]. Chen, H.; Chung,W.; Xu, J.J.;Wang, G.; Qin, Y.; Chau, M. Crime data mining: A general framework and some examples. Computer 2004, 37, 50–56.

[2]. Liu, H.; Brown, D.E. Criminal incident prediction using a point-pattern-based density model. Int. J. Forecast. 2003, 19, 603–622.

[3]. Nakaya, T.; Yano, K. Visualising crime clusters in a space-time cube: An exploratory data-analysis approach using spacetime kernel density estimation and scan statistics. Transactions 2010, 14, 223–239.

[4]. Chandra, B.; Gupta, M.; Gupta, M.P. A multivariate time series clustering approach for crime trends prediction. In Proceedings of the 2008 IEEE International Conference on Systems, Man and Cybernetics, Singapore, 12–15 October 2008; pp. 892– 896.

[5]. Gouvis, R.C. Schools as Generators of Crime: Routine Activities and the Sociology of Place; American University: Washington, DC, USA, 2002.

[6]. Kautt, P.M.; Dennis, W.R. Schools as criminal hot spots primary, secondary, and beyond. Crim. Justice Rev. 2007, 32, 339– 357.

[7]. Engstad, P.A. Environmental Opportunities and the Ecology of Crime. In Crime in Canadian Society; Silverman, R.A., Teevan, J.J., Eds.; Butterworths: Toronto, ON, Canada, 1975.

[8]. Brantingham, P.; Brantingham, P. Criminality of place. Eur. J. Crim. Pol. Res. 1995, 3, 5–26.

[9]. Block, R.L.; Davis, S. The Environs of Rapid Transit Stations: A Focus for Street Crime or Just Another Risky Place? In Preventing Mass Transit. Crime; Clarke, R.V., Ed.; Criminal Justice Press: Monsey, NY, USA, 1996; pp. 237–257.

[10]. Block, R.L.; Block, C.R. The Bronx and Chicago: Street Robbery in the Environs of Rapid Transit Stations. In Analyzing Crime Patterns: Frontiers in Practice; Goldsmith, V., McGuire, P.G., Mollenkopf, J.H., Eds.; Sage Publications: Thousand Oaks, CA, USA, 2000; pp. 137–152.

[11]. Salakhutdinov, R.; Mnih, A.; Hinton, G. Restricted Boltzmann machines for collaborative filtering. In 2007 Proceeding, Proceedings of the 24th International Conference on Machine Learning, Corvallis, OR, USA, 20 June 2007; Association for Computing Machinery: New York, NY, USA, 2007; pp. 791–798.

[12]. Lee, S.; Jung, S.; Lee, J. Prediction model based on an artificial neural network for user-based building energy consumption in South Korea. Energies 2019, 12, 608.

[13]. Huang, W.; Foo, S. Neural network modeling of salinity variation in Apalachicola River. Water Res. 2002, 36, 356–362.

[14]. Hwang, J.T. A study on target selection of burglars, robbers and thieves. Korean Inst. Criminol. 2004, 04-26, 217–247.

[15].Cohen, J. Statistical Power Analysis for the Behavioral Sciences, 2nd ed.; Lawrence Earlbaum Associates: Hillsdale, NJ, USA, 1988.

[16]. Stucky, T.D.; Ottensmann, J.R. Land use and violent crime. Criminology 2009, 47, 1223–1264.

[17]. Kim, D.K.; Yoon, Y.J.; Ahn, K.H. A study on urban crime in relation to land use patterns. J. Korea Plan. Assoc. 2007, 42, 155–168.
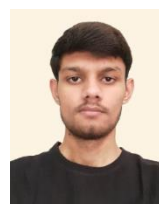
## AUTHORS

**Ajeet Singh**

He working as Assistant Professor in Department of Computer Science and Engineering at Moradabad Institute of technology, Moradabad, UP, India. He is pursuing his Phd from Faculty of Engineering and Technology at Rama University Uttar Pradesh, Kanpur (India) in Computer Science and Engineering. He did his B.tech and M.tech in Computer Science. His Research area is AI Algorithms and Deep Learning.

**Kunal Sharma**

My name is Kunal Sharma. Currently, I am pursuing my Bachelor of Technology in Computer Science and Engineering from Moradabad Institute of Technology. I am skilled in Python, C, HTML, CSS, and JavaScript. I have completed various projects related to these skills and my areas of interest are Machine Learning and Data Science.
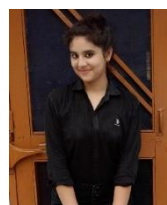
**Md Vadiyat Naqvi**

I am Md. Vadiyat Naqvi. A computer science engineering student at Moradabad Institute of Technology Affiliated to Abdul Kalam Technical University(AKTU) Lucknow. I am in 4th year and strong skilled in HTML, CSS, JavaScript, Java, C/C++, Python, MySQL, Responsive Design and worked on many projects related to these skills.

**Khushi Tyagi**

I am Khushi Tyagi. Currently, I am pursuing my bachelor of technology in computer science from Moradabad Institute of technology which is affiliated from Aktu, Lucknow. I am good in python programming language. Current I am working on Machine learning and deep learning. Apart from technical skills, I have a leadership and problem solving skills .

**Kumkum Joshi**

I am Kumkum Joshi. Currently, I am pursuing my bachelor of technology in computer science from Moradabad Institute of technology which is affiliated from Aktu, Lucknow. I am good in python programming language. Current I am working on Machine learning and deep learning.