# Edge Assisted Crime Prediction and Evaluation Framework for Machine Learning Algorithms

Apurba Adhikary[1], Saydul Akbar Murad[2], Md. Shirajum Munir[1], and Choong Seon Hong[1]*

[1]*Department of Computer Science and Engineering, Kyung Hee University,*
*Yongin-si 17104, Republic of Korea*
[2]*Faculty of Computing, University Malaysia Pahang, Pahang, Malaysia*
E-mail: apurba@khu.ac.kr, MCN21001@student.ump.edu.my, munir@khu.ac.kr, cshong@khu.ac.kr

*Abstract*—The growing global populations, particularly in major cities, have created new problems, notably in terms of public safety regulation and optimization. As a result, in this paper, a strategy is provided for predicting crime occurrences in a city based on historical events and demographic observation. In particular, this study proposes a crime prediction and evaluation framework for machine learning algorithms of the network edge. Thus, a complete analysis of four distinct sorts of crimes, such as murder, rapid trial, repression of women and children, and narcotics, validates the efficiency of the proposed framework. The complete study and implementation process have shown a visual representation of crime in various areas of country. The total work is completed by the selection, assessment, and implementation of the Machine Learning (ML) model, and finally, proposed the crime prediction. Criminal risk is predicted using classification models for a particular time interval and place. To anticipate occurrences, ML methods such as Decision Trees, Neural Networks, K-Nearest Neighbors, and Impact Learning are being utilized, and their performance is compared based on the data processing and modification used. A maximum accuracy of 81% is obtained for Decision Tree algorithm during the prediction of crime. The findings demonstrate that employing Machine Learning techniques aids in the prediction of criminal events, which has aided in the enhancement of public security.

*Index Terms*—Machine Learning, Edge Computing, Crime Prediction, Impact Learning, Decision Tree, KNN, MLP

## I. INTRODUCTION

One of the primary concerns of the global population is public safety. Several causes, such as the rapid pace of urbanization, have led to the rise in worry. The migration of people to cities has been well-known in recent years, and according to UN predictions, over 70% of the world's population will live in cities by 2050 [1]. In addition, according to the Global Terrorism Database, which defines terrorist attacks as "acts of violence by non-state actors perpetrated against civilian populations, intended to cause fear, in order to achieve a political objective," the number of terrorist attacks in the last decade was the highest ever recorded. Machine learning (ML) techniques are critical for smart city applications and may be used to reduce crime since they aid with difficulties concerning

urban development and the extraction of value from the data obtained [2].

This paper provides a graphical representation of crime in several areas of country, such as Bangladesh.. We used data from 2012 to 2019 to show the seniors. Using the crime prediction model, we observed that region 1 is more corrupted than other cities whereas region 2 is less corrupted. Using all of this information, we created a model that predicts crime in 2021 (check this). We present a crime comparison between 2019 and 2021. We have incidents to inform residents and municipal officials about the most dangerous places, therefore providing value to the community and increasing public safety. In this regard, this sort of prediction can be beneficial in a variety of ways, including more efficient and effective patrol route planning, as well as for tourists who are ignorant of the city's most hazardous locations.

Machine learning techniques such as Impact Learning [3], Decision Tree [4], K-Nearest Neighbors [5], and MLP Classifier are used to make the predictions. These algorithms are graded based on the data processing and transformation techniques are utilized. We summarize the key contribution as follows:

1) First, we have proposed a crime prediction and evaluation framework for ML algorithms of the network edge. That not only can learn historical analysis for crime detection but it also can evaluate based on the current data to protect the crime.
2) Second, we have implemented several ML algorithms, such as Decision Trees, Neural Networks, K-Nearest Neighbors, and Impact Learning on top of the proposed framework. In which, real crime data of the country are used to verify the effectiveness of the proposed framework by a comprehensive study.
3) Finally, our results show the efficacy to predict the crime, such as Murder, Speedy trial, Woman and Child Repression, and Narcotics in terms of accuracy to protect the society.

The rest of this paper is laid out as follows. The related work is described in Section II, and the Proposed Framework is presented in Section III. The results and discussions are then provided in Section IV, along with a description of the machine learning methods. Finally, we conclude in section V.
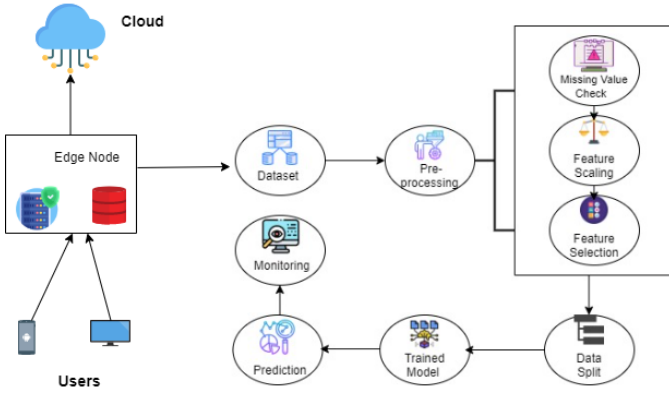
Fig. 1. A Crime Prediction and Evaluation Framework for Machine Learning Algorithms of Network Edge.

## II. RELATED WORK

A group of researchers used WEKA, an open-source data mining software, to compare violent crime patterns from the Communities and Crime Unnormalized Dataset provided by the University of California-Irvine repository with actual crime statistical data for the state of Mississippi obtained from neighborhoodscout.com [6]. Based on historical event data and demographic data, Luis et al. proposed a model for predicting crime occurrences in a city [7]. Another research group considered the creation of a crime prediction prototype model utilizing the decision tree (J48) algorithm. J48 algorithm detected the unknown category of crime data with 94.25287 percent accuracy based on the testing findings [8]. In another study, the authors employed machine learning and data science techniques to forecast crime in a Chicago crime data set [9]. The crime statistics were obtained from the Chicago Police Department's official website. In a different study, crime data from Vancouver for the previous 15 years was examined using two distinct data-processing techniques. When it comes to forecasting crime in Vancouver, the accuracy ranges from 39% to 44% [10]. Another research group proposed a crime prediction model based on communes (the regions or districts that make up the city of Buenos Aires) and utilized Python programming language for the prediction.Aside from [6–10], we are offering a comprehensive framework in this work whereby inputting the year, this model will provide the forecast of crime on the putted year.

## III. PROPOSED FRAMEWORK

We performed this assignment with the help of a machine learning system. To make use of the obtained data, we followed five steps [11], as shown in *Figure 1*. The high number of samples and enough varied variance help to preserve data quality. The proposed framework consists of six steps process. The first is data gathering, which is followed by monitoring. We have pre-processed the data after collecting it. In this step, we initially check for missing values before moving on to feature scaling. Finally, we choose the feature for label and feature data. We employ the ML technique after completing

TABLE I
DATASET ATTRIBUTES DETAILS INFORMATION.

| Attributes | Description | Type |
|---|---|---|
| Murder | Number of Murder in different city of the country. | Numerical |
| Speedy trial | Speedy trial in five Metropolitan areas of the country. | Numerical |
| Woman and Child Repression | Up and down of this crime is shown in number. | Numerical |
| Narcotics | Number of Narcotics in different city of the country. | Numerical |

the data split for training and testing. Based on the training data, this ML technique made a prediction. The proposed framework is physically deployed in the network edge.

### A. Edge Network

Edge networking is a distributed computing architecture that moves computation and data storage as close as feasible to the point of request to reduce latency and conserve bandwidth [12], [13]. Edge computing captures and processes data as close as feasible to the source of the data or intended event. It collects data using sensors, computer devices, and machines before sending it to edge servers or the cloud [14]. This data may be used to feed analytic and machine learning systems, provide automation capabilities, or provide visibility into the current condition of a device, system, or product [15], depending on the activity and desired consequence.

### B. Data Source

We gathered all of the information from the country's police website [16]. We've used dates ranging from 2012 to 2019. We identified a variety of statistics connected to crime here, but we chose some of the most relevant ones that are on the rise and highly prevalent in the country. *Table I* describe the collected data for our research purpose.

### C. Data Pre-processing

Data on crime prediction is pre-processed once various records are gathered [13]. Murder, Speedy trial, Woman Child Repression, and Narcotic are among the four observation values in this dataset for crime prediction. We got a lot of information from this website, but not all of it was wall-decorated. We used three distinct steps to pre-process the data for this.

- Missing Value check.
- Feature Scaling.
- Feature Selection.

**Missing Value check:** Values Management A missing value is incorrectly characterized in most situations as a value that was not saved in the example. In data, the lack of a value is a regular event. Furthermore, most foresight presentation methods are incapable of coping with any missing data. As a result, this problem should be addressed before modeling begins. We utilized the mean to replace a missing value. To find the mean, which is the same as the average value of a

data set, a calculation is needed. Subtract the total number of numbers from the total number of values in the data collection [15]. *Equation:1* is used to get the mean.

$$Mean = \frac{\sum_{i=0}^{n} x_n}{n} \tag{1}$$

**Feature Scaling:** Feature scaling or normalization is one of the most important procedures in machine learning approaches; without it, objective systems will not operate properly. Min-Max Scaling, Variance Scaling, Standardization, Mean Normalization, and Unit vectors are some of the feature scaling methods available. For this project, we used min-max normalization. The range in normalization in min-max in [0, 1] or [-1, 1] is given by *equation:2* and the fundamental formula of min-max in [0,1].

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)} \tag{2}$$

**Feature Selection:** Another job that must be completed prior to the deployment of a model is feature selection. The major objective of this process is to find a way to relate the feature to the target variables. We've removed a few highlights that aren't as important to the goal variable, while keeping the key highlights. The computational cost is reduced when the number of highlights is reduced. There are thirteen feature columns in our dataset.
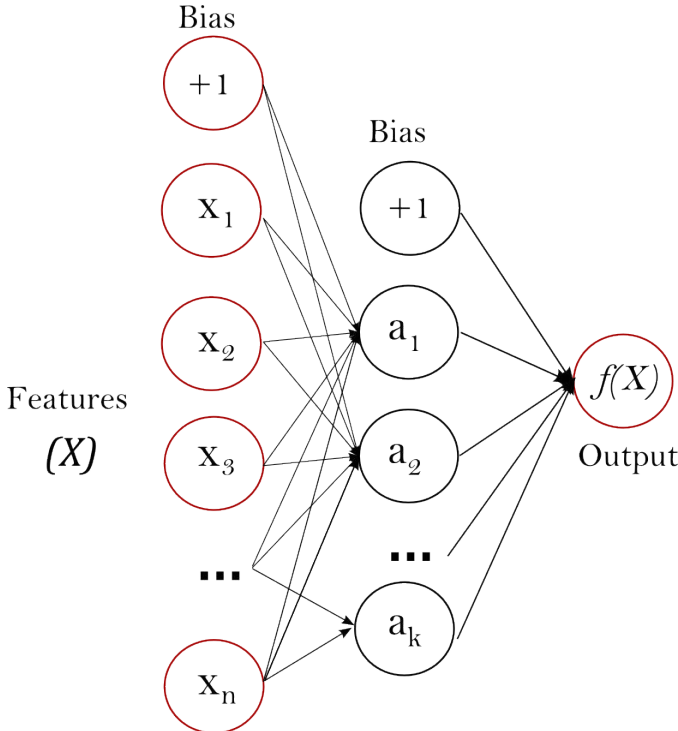


Fig. 2. Diagram of MLP Classifier [17].

### D. Data Split

Training data is the foundation for all machine learning algorithms. All of the information gathered has been split into two sections. The first is a training set, while the second

| Name of Algorithms | Description | Initial Parameters |
|---|---|---|
| K-Nearest Neighbors | The majority of votes from its fixed neighbors classify the new data point. | Nearest Neighbours: 5 |
| MLP Classifier | All input is processed through a series of hidden layers, and a final output is anticipated as a result. | Random state: 1 |
| Decision Tree | The more splits there are in a tree, the more information it captures about the data. | Class weight: None |
| Impact Learning | On the training dataset, a huge number of Epoch are run. Iteration is another name for Epoch. | Epoch: 2000 |

is a test set. We used 70% of the data for training and the remaining 30% for testing. This dataset is then trained using a machine learning model when this step is completed.

### E. Trained Model

We utilized four machine learning algorithms to find the most incredible accuracy. For this dataset, each method performs admirably. Based on the lowest error rate of the algorithms, the best-performing model is identified. *Table II* contains a summary of all implemented algorithms as well as their parameters.

**Impact Learning:** Impact learning is a knowledge learn approach that uses supervised classification and linear or polynomial regression. It also aids in the evaluation of rival data systems. It's rare to be able to understand the impact of independent features from a contest using this technique. In other words, the effects of the intrinsic rate of natural growth (RNI) [3] are taught. The RNI is represented by *equation:3* in this case.

$$\frac{dp}{dt} \approx rP \tag{3}$$

The impact learning equation is shown by *equation:4*:

$$Imp = (y' - (\frac{k\sum_{i=1}^{n} w_i x_i}{r - w_y k} + b))^{2/N} \tag{4}$$

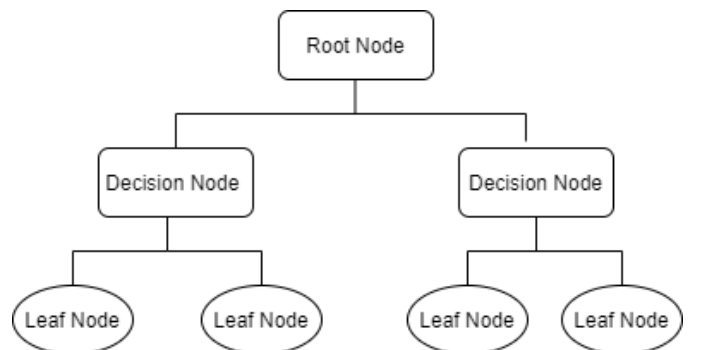**K- Nearest Neighbor Classifier:** The K-Nearest Neighbors



Fig. 3. Diagram of Decision Tree Classifier.

technique is one of the most fundamental Machine Learning algorithms, and it is based on the Supervised Learning methodology. The K-NN technique stores all available data and classifies new data points based on their similarity to previous data. This means that utilizing the K-NN method [5], new data may be swiftly classified into a well-defined category. The shortest distance between the selected neighbors is calculated using the KNN technique. The KNN uses the Euclidean distance function to calculate distances between existing data points and each new data point. *Equation:5* can be used to calculate Euclidean distance.

$$EuclideanDistance = \sqrt{\sum_{i=0}^{k}(x_i - y_i)^2} \qquad (5)$$

**MLP Classifier:** The term MLP Classifier refers to a Neural Network that uses Multi-layer Perceptron Classifiers. Unlike other classification methods such as Support Vectors or Naive Bayes Classifier, MLP Classifier uses an underlying Neural Network to perform classification. The perceptron is made up of two layers: an input layer and an output layer that are fully connected. In Figure 2 the first layer $X_1, X_2 \ldots X_n$ is the input layer and f(x) is the output layer. MLPs contain the same input and output layers, but they can have multiple hidden layers between them, as shown in *Figure 2*.

**Decision Tree:** It's a versatile forecasting technique that may be applied to a wide range of scenarios. In general, decision trees are an algorithmic technique for determining alternative ways to partition a data set depending on certain criteria. It's one of the most widely used methods of supervised learning [4]. The goal is to create a prototype that can learn and predict the value of a target variable using basic decision tree instructions. It is great for knowledge discovery since no parameter modifications are required. A Decision tree's two nodes are the Decision Node and the Leaf Node. Choice nodes are used to make any decision and contain multiple branches, whereas Leaf nodes are the consequence of those decisions and do not have any further branches. *Figure 3* describes the procedure of Decision Tree. It's divided into 3 steps. The first step is root node, and all others are sibling of this root node.

### F. Monitoring The Crime

Based on the predictions obtained from the prediction framework using ML algorithms of network edge [16], [17], responsible agency can monitor or evaluate the results.After monitoring the prediction results, required actions can be taken by the Government or police or the security agent.In addition, monitoring can be run time and emergency message can be broadcast via network to the acting agency i.e. government or police or the security agent.In this way, this proposed method will aid in the reduction of crime in a country.

### IV. RESULT AND DISCUSSION

The crime dataset was subjected to extensive testing in order to get the best output for crime prediction. To begin, the crime dataset is pre-processed in Google Colab, and a 30% portion
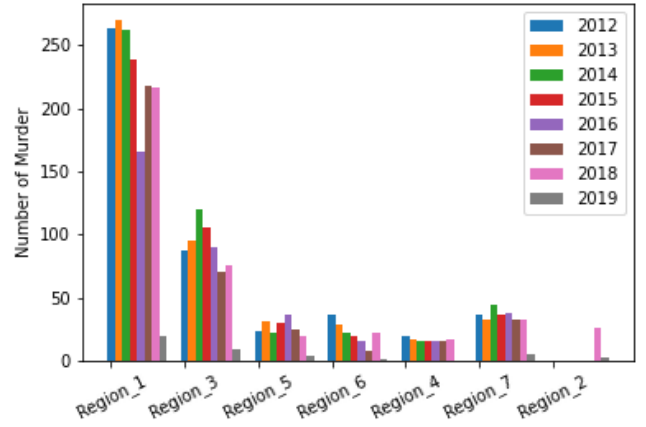
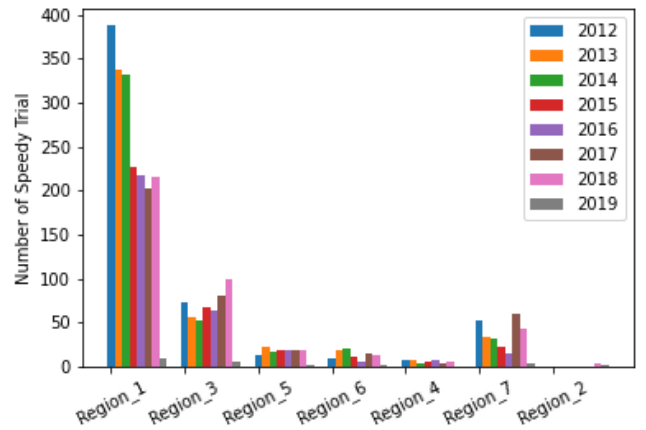

Fig. 4. The rate of Murder in various region.



Fig. 5. The rate of Speedy Trial in various region.

of the dataset is divided into training and test sets. We selected machine learning algorithms [18] and used the training data to create a classifier model for each algorithm which was used for testing. The results obtained indicate the performance of each classifier and the best classifier based on many metrics such as precision, accuracy, recall, and F-measurement for the given data set. Equation 6 is used to calculate accuracy.

$$Accuracy = \frac{TP + TN}{TP + TN.FP + FN} \qquad (6)$$

The algorithm's recall or sensitivity is represented as a percentage of all relevant results correctly categorized by the algorithm, which is expressed by the following *equation:7*.

$$Recall = \frac{TP}{TP + TN} \qquad (7)$$

Precision is the percentage of properly detected events or samples among those classified as positives. *Equation:8* demonstrates this.

$$Precision = \frac{TP}{TP + FP} \qquad (8)$$

The F1 score is calculated by taking the harmonic mean of accuracy and recall. *Equation:9* represents the F1 score.
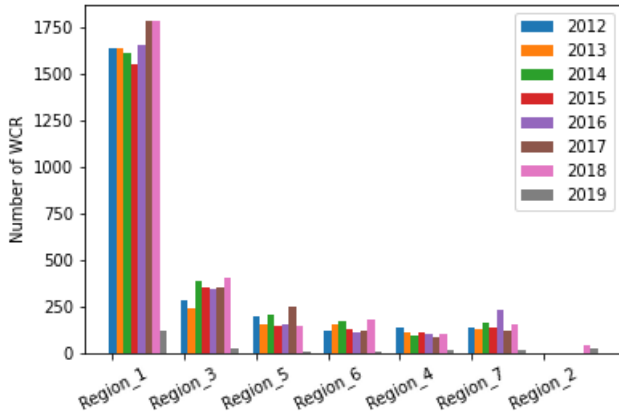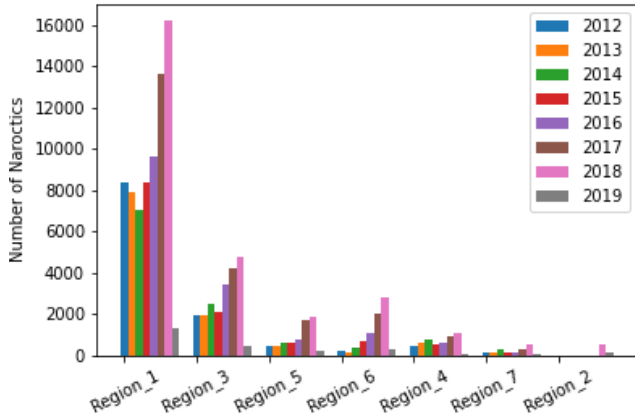
Fig. 6. The rate of WCR in various region.



Fig. 7. The rate of Narcotics in various region.

$$F1Score = \frac{2TP + TN}{2TP + FP + FN} \qquad (9)$$

True Positive (TP), False Positive (FP), True Negative (TN), and False Negative (FN) were considered for the evaluation of the performance metrics. The prediction skills of four artificial intelligence techniques were examined for the classification of crime prediction severity. *Table III* displays the classifier pre-
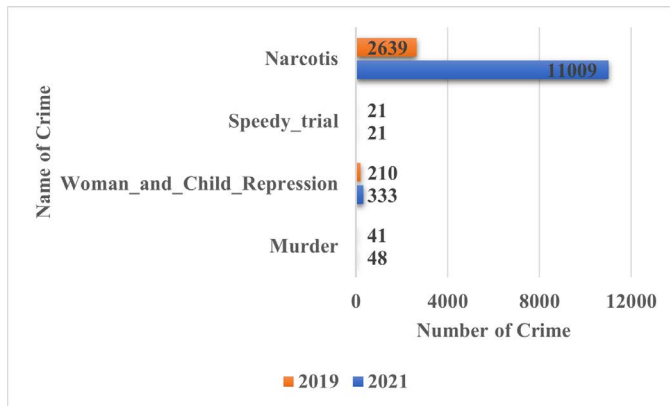


Fig. 8. Crime prediction in 2021 using Decision Tree and compression this crime what was happed 2019.

dictions. The prediction abilities of four artificial intelligence techniques were studied for crime severity categorization.

*Table III* compares the accuracy, precision, recall, and F1 score performance metrics of each method. We were able to obtain a high degree of accuracy for all of these algorithms by employing four robust methods. *Table III* shows that the decision tree has produced the highest accuracy of 81% whereas KNN classifier generated the lowest accuracy of 73%. In addition, decision tree algorithm produced the maximum value for all other performance metrics i.e. 73% of F1 score and Recall, and 78% of Precision. On the contrary, KNN Classifier algorithm produced the lowest value for the performance metrics i.e. 73% of accuracy, 66% of F1 Score, 69% of Recall and 70% of Precision. We also obtained good performance for all performance metrics for the other two algorithms, MLP classifier and impact learning. MLP classifier accuracy is 77%, while impact learning accuracy is 76%, which is comparable.The MLP classifier performed well in terms of F1 score and recall when compared to impact learning, but impact learning performed better in terms of Precision.

*Figure 4* shows that the murder rate at region 1 is quite high. The murder rate is lower at region 2.Region 3 has the second highest murder rate in the country.For another four Metropolitan areas, the rate is likewise in the average phase. The good news is that the murder rate in the metropolitan area is decreasing day by day as time passes.

We may also view a graph like this for a Speedy Trail in *Figure 5*. In comparison to other metropolitan areas, the rate in region 1 is much greater. The second post is for the region of the region 3. The rate of Speedy trials is lower in region 4 and region 2. The rate is in the intermediate phase for the remaining three metropolitan areas. However, the number of speedy trials is diminishing year after year. The total number of fast trials was 544 in 2012, but just 21 in 2019.

*Figure 6* shows that the WCR (Woman and Child Repression) rate was greater in 2018 than in 2019. This is a scenario that applies to all metropolitan areas. At region 1, we've always found the crime rate to be higher. This scenario is also applicable to WCR. Region 3 holds the second highest crime rate of all time. The rate is also greater in the region 5. For the next four locations, the rate is reduced. The most frustrating aspect is that, as time passes, the scenario of WCR remains unchanged, but in 2019, this ratio has decreased dramatically.

When we look at the other three crimes, we notice that the number of narcotics is significantly larger. Despite the fact that the rate was excessively high in 2018, it is still

in jeopardy.The most surprising aspect is that the number of Narcotics is increasing over time, which is a complete opposite scenario when compared to other types of crime. All locations have a greater rate of narcotics. As region 1 is a more densely inhabited location, the rate is greater than in other cities. *Figure 7* shows that narcotics are also in danger in region 3.However, we discovered a good scenario for Region 2 where the number of Narcotics is very low in comparison to other regions.

*Figure 8* depicts what the crime situation in the country will be in 2021. When we compare 2021 to 2019, we can see that the number of murders will increase somewhat in 2021. The number was 41 in 2019, but according to our estimate, it will be 48 by the end of the year. We obtained the same result in 2019 and 2021 for rapid trial, and the number is 21. We discovered a dreadful result for Narcotic. This graph shows that Narcotic will reach 11009 by the end of 2021, which is much higher than 2019. This finding might cause widespread concern among countrymen. In addition, we revealed that repression of women and children will rise in 2021. Our model predicts that by the end of 2021, the total number of crimes will have increased to its highest level. We looked at four different forms of crime and discovered that three of them will grow in 2021. In their [19] article, they only look at one crime in the country, but in our research, we look at four. Finally, we may conclude that this study will aid in raising environmental awareness.

## V. CONCLUSION

In this work, we have introduced a crime prediction and evaluation framework for machine learning algorithms of network edge. We collected data from 2012 to 2019 to analyze and evaluate our forecast. We used machine learning approaches to anticipate crime events, which can make a significant contribution to improving city public safety, which is a big problem in many cities across the world. It was fascinating to see how pre-processing, and transformation may affect the model's output, particularly when breaking the day into many time periods. Due to the provenance of the data, this solution was created for a specific city in the country. However, if equivalent data is made accessible, the technique may be applied to other cities. Based on the training set input for the four algorithms, we find the Decision tree method to be extremely successful and accurate in predicting crime data. The Decision Stump algorithm's poor performance could be attributed to a certain amount of randomness in the various crimes and associated features (shows a low correlation coefficient among the four algorithms); the KNN's branches are more rigid and only give accurate results if the test set follows the pattern modelled.

## REFERENCES

[1] "68% of the world population projected to live in urban areas by 2050, says UN — UN DESA — United Nations Department of Economic and Social Affairs." https://www.un.org/development/desa/en/news/population/2018-revision-of-world-urbanization-prospects.html (accessed Oct. 14, 2021).

[2] Y. Wu, W. Zhang, J. Shen, Z. Mo, and Y. Peng, "Smart city with Chinese characteristics against the background of big data: Idea, action and risk," Journal of Cleaner Production, vol. 173, pp. 60–66, Feb. 2018.

[3] M. Kowsher, A. Tahabilder, and S. A. Murad, "Impact-learning: A robust machine learning algorithm," in ACM International Conference Proceeding Series, pp. 9–13, Jul. 2020.

[4] Priyanka and D. Kumar, "Decision tree classifier: A detailed survey," International Journal of Information and Decision Sciences, vol. 12, no. 3, pp. 246–269, 2020.

[5] L. Jiang, Z. Cai, D. Wang, and S. Jiang, "Survey of improving K-nearest-neighbor for classification," Proceedings - Fourth International Conference on Fuzzy Systems and Knowledge Discovery, FSKD 2007, vol. 1, pp. 679–683, 2007.

[6] L. McClendon and N. Meghanathan, "Using Machine Learning Algorithms to Analyze Crime Data," Machine Learning and Applications: An International Journal, vol. 2, no. 1, pp. 1–12, Mar. 2015.

[7] Fonseca, Luis, F. C. Pinto, and S. Sargento. "An Application for Risk of Crime Prediction Using Machine Learning." International Journal of Computer and Systems Engineering 15.2, pp. 166-174, 2021.

[8] E. Ahishakiye, D. Taremwa, E. O. Omulo, and I. Niyonzima, "Crime Prediction Using Decision Tree (J48) Classification Algorithm," 2017. [Online]. Available: www.ijcit.com188

[9] S. K. Senthil Kumar, G. Adarsh, J. Shashank, and A. Sameer, "CRIME PREDICTION AND ANALYSIS USING MACHINE LEARNING", [Online]. Available: http://ijte.uk/

[10] S. Kim, P. Joshi, P. S. Kalsi, and P. Taheri, "Crime Analysis Through Machine Learning," in 2018 IEEE 9th Annual Information Technology, Electronics and Mobile Communication Conference, IEMCON 2018, pp. 415–420, Jan. 2019.

[11] S. A. Murad, Z. R. M. Azmi, Z. H. Hakami, N. J. Prottasha, and M. Kowsher, "Computer-aided system for extending the performance of diabetes analysis and prediction," pp. 465–470, Sep. 2021.

[12] L. U. Khan, I. Yaqoob, N. H. Tran, S. M. A. Kazmi, T. N. Dang, C. S. Hong, "Edge Computing Enabled Smart Cities: A Comprehensive Survey," IEEE Internet of Things Journal, Vol.7, Issue 10, pp.10200-10232, Oct. 2020.

[13] M. S. Munir, S. F. Abedin and C. S. Hong, "Artificial Intelligence-based Service Aggregation for Mobile-Agent in Edge Computing," 2019 20th Asia-Pacific Network Operations and Management Symposium (APNOMS), pp. 1-6, 2019.

[14] M. S. Munir, S. F. Abedin, N. H. Tran and C. S. Hong, "When Edge Computing Meets Microgrid: A Deep Reinforcement Learning Approach," in IEEE Internet of Things Journal, vol. 6, no. 5, pp. 7360-7374, Oct. 2019.

[15] M. S. Munir, S. F. Abedin, D. H. Kim, N. H. Tran, Z. Han and C. S. Hong, "A Multi-Agent System toward the Green Edge Computing with Microgrid," 2019 IEEE Global Communications Conference (GLOBE-COM), pp. 1-7, 2019.

[16] "Bangladesh Police". https://www.police.gov.bd/en/crime statistic/year/2019 (accessed Oct. 13, 2021).

[17] "sklearn.tree.DecisionTreeClassifier-scikit-learn 1.0 documentation."https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html (accessed Oct. 13, 2021).

[18] A. J. M. Muzahid, S. F. Kamarulzaman, and M. A. Rahman, "Comparison of PPO and SAC Algorithms Towards Decision Making Strategies for Collision Avoidance Among Multiple Autonomous Vehicles," pp. 200–205, Sep. 2021.

[19] M. A. Awal, J. Rabbi, S. I. Hossain, and M. M. A. Hashem, "Using linear regression to forecast future trends in crime of Bangladesh." 2016 5th International Conference on Informatics, Electronics and Vision (ICIEV). IEEE, 2016.