

## CRIME PREDICTION USING K-NEAREST NEIGHBOURING ALGORITHM

N. Naga Swathi\*<sup>1</sup>, Sarah Vineela Cole\*<sup>2</sup>, Polani Nikhil Manikanta\*<sup>3</sup>,

Shaik Ashmiya Jafreen\*<sup>4</sup>, Shaik Salma\*<sup>5</sup>, Talari Manoj\*<sup>6</sup>

\*<sup>1</sup>Asst., Professor, Dept. Of Electronics And Communication Engineering, Bapatla Engineering College, Bapatla, India.

\*<sup>2,3,4,5,6</sup>Student, Department Of Electronics And Communication Engineering, Acharya Nagarjuna University, Bapatla Engineering College, Bapatla, India.

### ABSTRACT

Violations are a social bothering and taken a toll our society deeply in a few ways. Any inquire about that can offer assistance in tackling crimes rapidly will pay for itself. Almost 10% of the offenders commit almost 50% of the violations [9]. The framework is prepared by feeding past a long-time record of violations taken from legitimate online entry of India posting different violations such as kill, kidnapping and kidnapping, dacoits, burglary, burglary, assault, and other such violations. As per information of Indian insights, which gives data of different crime of past 14 a long time (2001-2014) a relapse model is made and the crime rate for the taking after a long time in various states can be anticipated [8]. We have utilized supervised, semi-supervised and unsupervised learning procedure [4] on the crime records for information revelation and to assist in expanding the prescient precision of the wrongdoing. This work will be supportive to the nearby police stations in crime concealment.

**Keywords:** Crime Prediction, Machine Learning, KNN, Regression, Machine Learning, Styling.

### I. INTRODUCTION

Criminal movement is slowly rising in India and includes a significant and negative social affect [3]. The later spurt in the country has put everyone pondering as to what will happen within the future. Cases of kill, snatching, assault, and fatal mishaps have skyrocketed. The require of the hour is to make individuals of the country realize the issue. Machine learning headways and profound learning calculations can find modern designs in different information sets and uncover modern information. Wrongdoing forecast and distinguishing offenders are the one of the beat need issues to the police office because there's a colossal sum of information related to crime that exists. There's a require for innovation through which the case-solving can be quicker [3]. The thought behind this venture is that wrongdoings can be effectively anticipated once we are able to sort through a gigantic sum of information to discover patterns that are valuable to arranging what is required [1]. The later advancements in machine learning makes this task conceivable. One will donate date, time, area (longitude, latitude) as input and the yield will be created which can give us data around which wrongdoing is likely to happen in that region. It essentially gives us the hotspots of wrongdoing [5]. The data is taken considering the time and sort of wrongdoing that happened within the past. KNN calculation at that point employments its approach which accept that comparative things exist in near nearness and classifies modern cases based on likeness measures.

Classes of Crimes are:

- Act 363 - Kidnapping
- Act 379 - Robbery
- Act 13 - Gambling
- Act 302 - Murder
- Act 279 - Accident
- Act 323 - Violence

This prediction, on the off chance that put to great utilize, is of awesome offer assistance in investigating cases that have happened. It can be utilized to suppress the wrongdoings by introducing a few measures on the off chance that we know what sort of wrongdoing is planning to happen already. This will indirectly offer assistance decrease the rates of wrongdoings and can offer assistance to improve security in such required ranges [2].

## II. RELATED WORK

It is seen that numerous of machine learning models are built on datasets of diverse cities having diverse special highlights, so suspicion is diverse in all cases. Classification models have been executed on different other applications like forecast of climate, in managing an account, funds additionally in security [3].

In [4] recognizable proof of offenders by utilizing classification methods and wrongdoing expectation was done utilizing information set of six cities of Tamil Naduby utilizing KNN classification, K-Means clustering, Agglomerative progressive clustering, and DBSCAN clustering calculations. In [5], they utilized a show whose fundamental point was to utilize a dataset where the information positions were isolated into different classes to induce clarity of a modern test positions. Utilizing highlights like Day, Date, Year of the wrongdoing utilizing KNN - calculation it is found to be 40% exact. Their demonstrate utilized strategies like Calculated Regression, Choice Trees, Bayesian Strategies and Bolster Vector Machine [9]. Python came into utilize for preparing information, make relapse investigation and conclude the categories for test dataset, to urge the most excellent relationship between the highlights (Hour, Longitude, Scope, Day of the Week, Week, Month) and the goal esteem (Divisions of Wrongdoing). All-important values were changed into parallel values by making the values of the highlights values into partitioned unused traits and convert values into either a or 1. There were numerous trails of different Relapse strategies were utilized on the preparing dataset by partit into two sets; preparing and testing, both validation and cross-testing were conducted, the strategy with the least misfortune was connected to induce the comes about for the test information.

## III. PROPOSED WORK

### A. Processing data:

At first, information was pre-processed by expelling all invalid values and columns that are pointless [2]. The dataset that was used could be a adjustment of the first dataset that was obtained by scratching the police site of a city Indore in Madhya Pradesh. It was handled different times and they dropped highlights such as police station, station number, Complainant title & address, blamed title & address. There were minor alterations made in their last dataset. importance of highlights was calculated by the Additional Trees Classifier work which made a difference us in neglecting the unnecessary qualities (allude Table 1). Extra Trees Classifier may be a sort of outfit learning technique which takes a indicated sum of information (esteem of n-estimators) and calculates the significance of each and every highlight independently. Figure 1 appears the significance of each highlight in through a bar graph.

**Table 1: Importance of Features**

Features	Importance
Hour	0.29611912
Latitude	0.361180
Longitude	0.2677035
Year	0.00000
Month	0.0013
Week of the year	0.00442

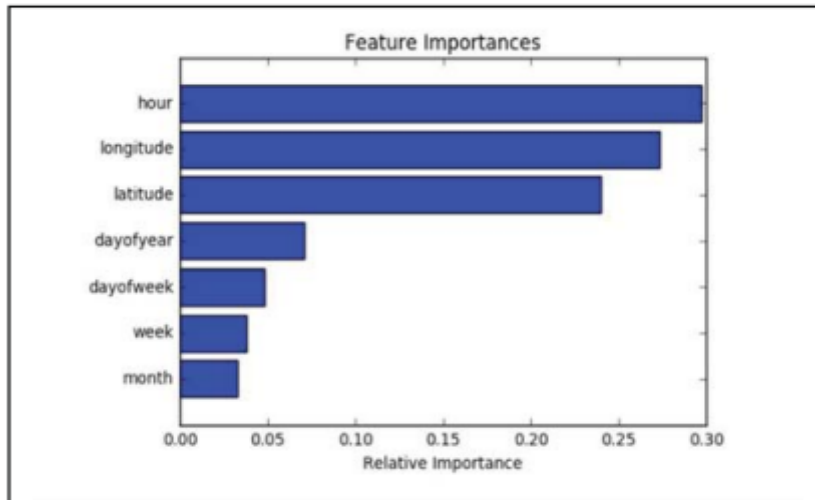


Figure 1: Importance of Features

	hour	dayofyear	act379	act13	act279	act323	act363	act302	latitude	longitude
0	21.0	59.0	1	0	0	0	0	0	22.737260	75.875987
1	21.0	59.0	1	0	0	0	0	0	22.720992	75.876083
2	10.0	59.0	0	0	1	0	0	0	22.736676	75.883168
3	10.0	59.0	0	0	1	0	0	0	22.746527	75.887139
4	10.0	59.0	0	0	1	0	0	0	22.769531	75.888772

Figure 2: Final data set

The traits of slightest significance were dropped (year, month, week of the year, etc.). The ultimate dataset (allude Figure 2) presently has four qualities with hour, day of the year, longitude and scope of the city. After the ultimate dataset was made, another sub-set (allude Figure 3) was made by utilizing SQL (allude Figure 4). A SNS heatmap was created (allude Figure 4) to urge a harsh thought approximately how the wrongdoings are varying with regard to days of a month. A warm outline is like data investigation computer program which takes the assistance of colours the way comparable to a bar chart which employments stature and width as a data visualization apparatus. [8] Our perceptions were that act 323 i.e., savagery had happened most towards the month-end whereas act 279 i.e., mischances happened then again.

day	act	frequency
0	1 act379	121
1	1 act13	22
2	1 act279	88
3	1 act323	66
4	1 act363	33
5	1 act302	0
6	3 act379	66
7	3 act13	0
8	3 act279	121
9	3 act323	66

Figure 3: Sub Data set

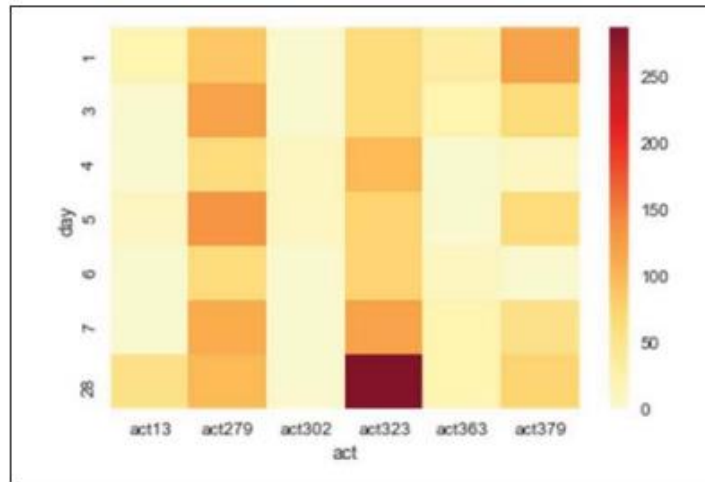


Figure 4: SNS Heatmap

**B. KNN Algorithm:**

The another step was to choose which calculation to utilize. K Nearest Neighbour Classifier could be a directed machine learning calculation valuable for classification issues. It works by finding the separations between a inquiry and all the examples within the data, selecting the desired illustrations that are closest to the inquiry, and after that votes for the foremost visit label. It isn't parametric which infers that it does not make any supposition on the essential information distribution. To put it in basic words, the model structure is chosen by the data. It's pretty valuable since in reality, most of the information does not take after the normal hypothetical standards made [4]. Hence, we chosen to utilize K-Nearest-Neighbour Calculation

**C. Applying Algorithm:**

The trained and test information were part within the taking after way. The after pie chart(refer Figure 5) demonstrates:

- Preparing dataset comprises of 80% information.
- Testing dataset comprises of 20% information.

This test set serves as a intermediary for unused information. One ought to make beyond any doubt that it is agent of the information set as a entire. To predict the esteem of k, a chart was plotted by utilizing the Elbow strategy. [7]

Our Test set serves as a method

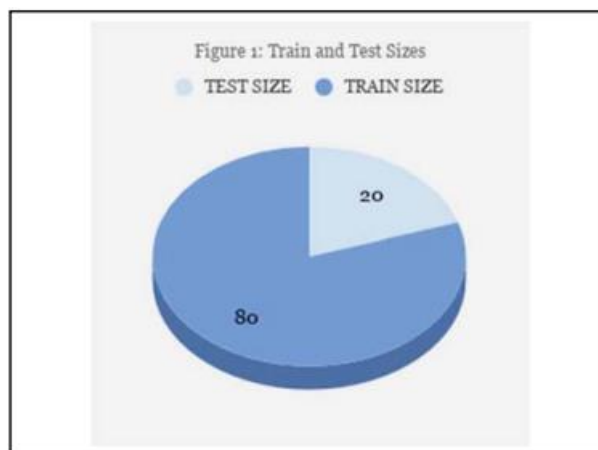


Figure 5: Split Sizes

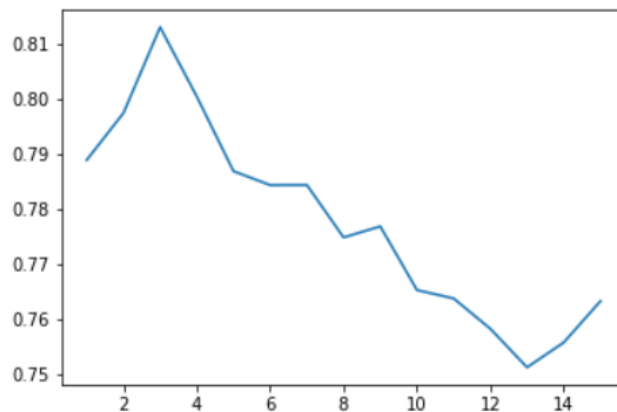
Strategy is exceptionally broadly utilized strategy which makes a difference in deciding the ideal esteem of k. The elbow strategy runs k-means clustering on the dataset for a assortment of values for k (e.g. 1-15) and after that for each esteem of k, it works out an normal score for all clusters. After analysing the chart, the extend in which the blunder rate was least came out to be 1-15. Besides, we checked all its values in extend 1-15, but

from the values of k extending 1-13, the precision remained consistent. Thus, we chosen k=3 that had a place from the run 1-13. To calculate MAE and RMSE values, we required test, train, and predicted values of both x and y.

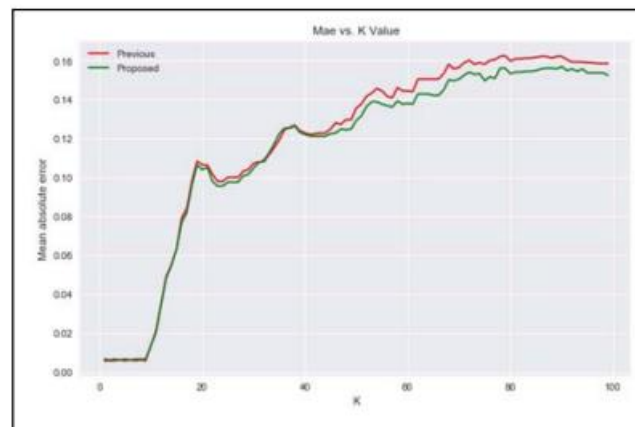
RMSE is a quadratic scoring rule which figures the average magnitude of the error. It is the square root of the square differences calculated between expectation and actual observation. MAE and RMSE both express average model prediction error in units of the interest variable. Provided that the errors are squared until they are averaged, large errors are assigned a fairly high weight. RMSE avoids making use of absolute value, which is unwanted in many mathematical calculations. After calculating both the values, we plotted two graphs.

$$MAE = \frac{1}{n} \sum_{j=1}^n |y - \hat{y}|$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^N (Predicted_i - Actual_i)^2}{N}}$$



**Figure 6: Mean Classifier**



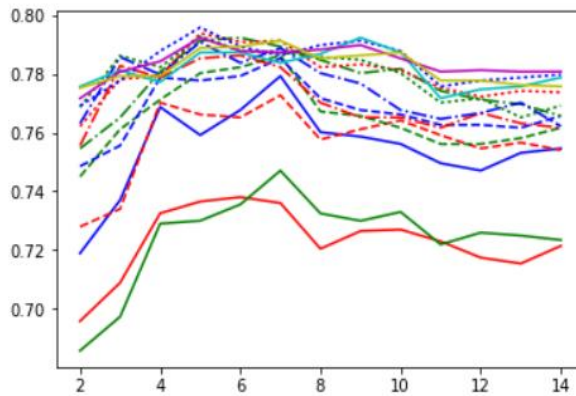
**Figure 7: MAE vs K**

The chart (allude Figure 6) shows the Cruel Supreme Error (Y-axis) and values of k (X-axis). MAE measures the average sum of the mistakes in a set of forecasts, without bearing in intellect the course. Second chart (allude Figure 7) appears Root Cruel Square Error (X-axis) and values of k (Y-axis). The Root Cruel Square Mistake could be a commonly utilized degree of the contrasts between the anticipated values by a show and the watched values.

The ruddy bend shows past work and the green bend indicates proposed work. Clearly, the cruel supreme blunder and root cruel square mistake is diminished when compared with previous work. The ruddy bend demonstrates past work and the green bend shows proposed work.

RMSE is diminished when compared with past work. Finally, after calculating all blunder and cruel values, the accuracy score of the program was calculated. To calculate precision for persistent factors, Cruel Supreme Mistake (MAE) and Root cruel squared mistake (RMSE) are two of the most common measurements that are utilized.

**IV. RESULTS AND DISCUSSIONS**



**Figure 8:** Error Rate vs K value

**Table 2:** Results

	Previous	Proposed
Mean Absolute Error	0.1598	0.0167635
Accuracy	0.5834	83.592574
KNN score	0.9323	0.9951

K = 12 made a difference procure the most elevated precision as conceivable. The above chart (allude Figure 8) shows the RMSE esteem plotted against K esteem. The yellow bend is for proposed work and blue is for past work. An increment in k-value results in expanded root cruel square mistake. Subsequently the esteem of k was picked from run 1-15 since that's the as it were range with least blunder. The past work had included extra variables which did not appear vital in our case. Comparison between the comes about from previous work and proposed work is shown in Table 2

Because it is evident from the chart, presently one can certainly say that the blunder was diminished in this way expanding the exactness of the program [3]. The designs of wrongdoing are not same each time designs continuously changes after time to time. The framework was prepared to memorize utilizing some particular inputs. So, the method by itself learns diverse changes that come within the pattern of wrongdoing after analysing them. Too, we cannot ignore the reality that wrongdoing components alter with time [3].

**V. CONCLUSION**

This inquires about work offers a way to anticipate and foresee crimes and fakes inside a city. It centers on having a wrongdoing prediction apparatus that can be accommodating to law requirement. This paper is pointed at expanding the expectation exactness as much as conceivable. As compared to the past work, this work was fruitful in accomplishing the most noteworthy precision in prediction. The values of RMSE and MAE were decreased significantly. Along the way, numerous designs of criminal activities in different regions which can be accommodating for criminal investigation were known. This design has much more prominent importance than we realize. The KNN framework makes a difference law implementing offices for made strides and correct wrongdoing analysis. By navigating through the wrongdoing dataset, we have to discover out distinctive reasons that lead to wrongdoing. Since this paper is bearing in intellect as it were a few restricted variables, full accuracy cannot be fulfilled.

For getting more exact results in expectation we got to discover out more wrongdoing attributes of places rather than setting as it were certain qualities. Thus, distant this framework was prepared utilizing certain traits, but we can take account of more components to progress precision. In the future, this work can be

extended to have created classification calculations to detect criminals more efficiently. The wrongdoing rates that are expanding continuous may go down within the future due to such expectation procedures.

## VI. REFERENCES

- [1] Kim, Suhong, Param Joshi, Parminder Singh Kalsi, and Pooya Taheri. "Crime Analysis Through Machine Learning." In 2018 IEEE 9th Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON), pp. 415-420. IEEE, 2018.
- [2] Shah, Riya Rahul. "Crime Prediction Using Machine Learning." (2003).
- [3] Lin, Ying-Lung, Tenge-Yang Chen, and Liang-Chih Yu. "Using machine learning to assist crime prevention." In 2017 6th IIAI International Congress on Advanced Applied Informatics (IIAI-AAI), pp. 1029-1030. IEEE, 2017.
- [4] M. V. Barnadas, Machine learning applied to crime prediction, Thesis, Universitat Politècnica de Catalunya, Barcelona, Spain, Sep. 2016.
- [5] Crime Prediction Using Machine Learning Sacramento State [athena.ecs.csus.edu > ~shahr > progress\\_report](http://athena.ecs.csus.edu/~shahr/progress_report) by RR Shah - 2003.
- [6] Williams, Matthew L., Pete Burnap, and Luke Sloan. "Crime sensing with big data: The affordances and limitations of using open-source communications to estimate crime patterns." *The British Journal of Criminology* 57, no. 2 (2017): 320-340.
- [7] Agarwal, Shubham, Lavish Yadav, and Manish K. Thakur. "Crime Prediction Based on Statistical Models." In 2018 Eleventh International Conference on Contemporary Computing (IC3), pp. 1-3. IEEE, 2018.
- [8] Nakaya, Tomoki, and Keiji Yano. "Visualising crime clusters in a spacetime cube: An exploratory data analysis approach using space time kernel density estimation and scan statistics." *Transactions in GIS* 14, no. 3 (2010): 223-239.