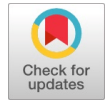


Crime Analytics using Machine Learning

Suman Acharya



Abstract: Crime is one of the most significant and pervasive problems in our society, and preventing it is a crucial duty. A large number of crimes are perpetrated each day. Maintaining and analyzing crime data to forecast and solve crimes is the current issue. This project analyzes a large dataset of crimes and predicts future crimes based on conditions. This project uses data science and machine learning for India's crime data prediction. Thus, Decision Tree, Logistic Regression, Multi-Regression, k-NN, Lasso & Ridge, and Random Forest are all involved in the supervised classification problem. Predicting crimes and classifying effective pattern detection and visualization equipment Utilizing crime data trends from the past allows us to correlate aspects that may help us comprehend the breadth of crimes in the future. This study uses visualization and machine learning methods to estimate future crime rates. First, raw datasets were processed and displayed

Keywords: Predictive Analytics, Crime Analytics, Machine Learning, Performance Enhancement, Pattern detection, Decision Learning.

I. INTRODUCTION

Crime poses the greatest threat to humanity. There are numerous crimes that occur at regular intervals. Perhaps it is growing and spreading rapidly and widely. From small villages and towns to major metropolitan areas, criminal activity is prevalent. There are various types of crimes, including robbery, murder, rape, assault, imprisonment, and kidnapping, etc. [1][2]. Crime prediction helps analysts visualize criminal networks, reduce risks, and boost productivity. [3][4]. A good prediction technique helps predict crime rates, evolve crime data sets faster, and track crime analysis resources. Crime analysis can be done using machine learning techniques. Machine learning approaches use computers and mathematics to program systems to operate. These methods help prevent and identify crime. Crime analysis includes pattern extraction, prediction, and detection. With so much crime data, the police force struggles to predict crime. Technology is required for faster case resolution. train a prediction model. The test dataset will validate the training dataset. Based on accuracy, smarter algorithms will build the model [5][6]. This work helps Indian law enforcement organizations predict and detect crimes more accurately, lowering crime rates. Machine learning algorithms have greatly improved crime prediction based on prior data.

Manuscript received on 07 January 2023 | Revised Manuscript received on 14 March 2023 | Manuscript Accepted on 15 March 2023 | Manuscript published on 30 March 2023.

*Correspondence Author(s)

Prof. Suman Acharya*, Assistant Professor, University of Engineering and Management, Jaipur (Rajasthan), India. E-mail id: Suman.acharya@uem.edu.in, ORCID ID: <https://orcid.org/0009-0002-2039-7239>

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an [open access](https://creativecommons.org/licenses/by-nc-nd/4.0/) article under the CC-BY-NC-ND license <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

The goal of this project is to use machine learning models to analyze and predict crimes over the next few months. It creates a model to estimate monthly crime by type. In this research, a number of machine learning models, such as regression techniques, K-NN, and boosted decision trees, will be utilized to predict criminal behavior. A month-wise crime analysis can be performed to comprehend the crime pattern [7]. Using various visualization tools and graphs can aid law enforcement organizations in detecting and predicting crimes with greater precision. This will indirectly aid in reducing crime rates and also enhance security in such necessary places. We utilized a dataset of crime reports for India from 2022 that was made available on the official website for datasets. The data, which originated from the Indian Police Department, had over 5473 records, or data points [7]. Each data point consists of 14 attributes representing diverse information on the major head, minor head, crimes committed in the past years and months, and the severity of the reported offense [7].

II. LITERATURE SURVEY

Machine learning models forecast crime using India's crime data collection. This paper compares KNN, SVM, and regression models. Dataset and feature selection affect prediction. KNN predicts 80%, Linear Regression 91%, SVC 86%, Lasso and Ridge Regression 85%, Logistic Regression 87%, and SVC 84%. [8][9][10]. Crime pattern detection is a big issue. Understanding datasets is also crucial. To avoid wasting resources on erroneous or missing numbers, we took the mean. Additionally, employing a strategy for categorizing the crime rate as high, medium, or low No one has identified the types of crimes that can occur or their likelihood of occurring. Crime analysis and prediction are crucial activities that can be optimized through the use of a variety of approaches and processes. Numerous researchers conduct extensive studies in this field. Existing work is confined to identifying the number of crimes committed in the current month using datasets [11][12] [13].

III. PROPOSED SYSTEM

The proposed system employs the Lasso-KNN combined technique to get greater precision and outcomes than other superior algorithms such as random forest, SVR, KNN, and decision tree classifiers. Lasso employed two-thirds of the Crime dataset for training, and then we trained Lasso's projected value again with KNN to achieve more accurate results. a collection of dataset-record decision trees. Lasso and KNN have a precision of 85%, random forest 84%, and SVR 90%. Our primary goal is to increase the precision of the KNN regressor method from 80% to 85%, predict accurate results for our Crime dataset using Lasso-KNN, and avoid overfitting and underfitting.

Our proposed approach employs Lasso and KNN, which are more accurate than previous regression algorithms by 85%. The next step uses Lasso and SVR, which are 91% more accurate than the previous regression algorithms and Lasso-KNN. Numerous graphs, such as bar graphs, may be applied in this representation. Using the matplotlib package from Sklearn [14]. The crime dataset is analyzed through the use of graphs

IV. METHODOLOGY

- **Data Collection Methods:**

Using important data sources, data collection was done for crime prediction. These are:

- Social media analysis
- Records of crime from the website of the Indian government
- Social media analysis
- newspapers;
- CRDs (Call Data Records)

The Indian Police Department's "Crime Data 2022" is the largest data set kept by a law enforcement agency. A sizable number of crime prediction systems are being tried in India. Social media is another important method of gathering information for crime prediction [7].

To accurately estimate the victims of the current month, we applied certain machine learning techniques. Using a combination of KNN and Lasso also improves the accuracy of KNN. SVR and the Random Forest Regressor are used to make accurate predictions. We estimate the appropriate number of victims for the months of August and December of the current year using databases of all crimes committed in those months.

- **Approach:**

Lasso-KNN:

Here, we utilize a combination approach of Lasso and KNN to improve the accuracy of KNN from 80% to 85%. Specifically, we first train the model with Lasso and then train it again with KNN in order to increase the accuracy to 85%.

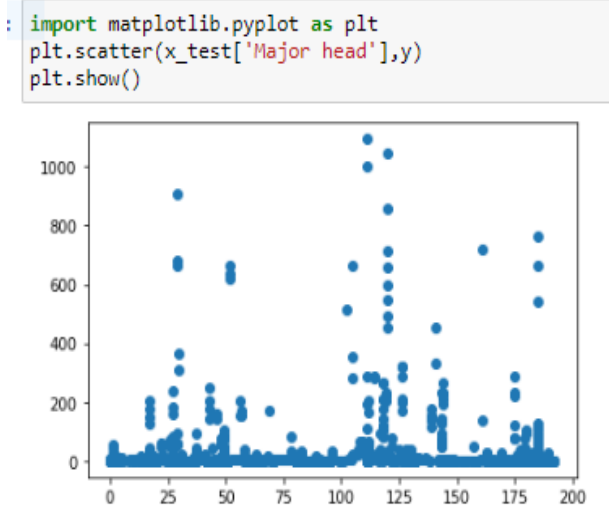


Fig 4.1: Scatter plot of Lasso-KNN of Crime Dataset Between major Head and Y Predict.

Lasso-SVR:

Here, we use a combination of Lasso and SVR to make the previous proposed method more accurate, from 90% to 91%. To be more specific, we train the model with Lasso first and then train it again with SVR to get the accuracy up to 91%.

```
In [48]: import matplotlib.pyplot as plt
plt.scatter(x_test['Major head'],y)
plt.show()
```

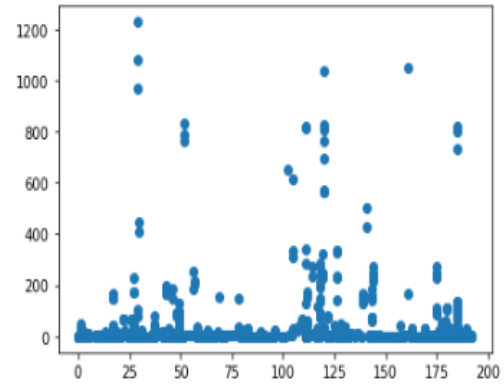


Fig 4.2: Scatter Plot of Lasso-SVR of Crime Dataset Between Major Head and Y Predict.

In this project, we viewed the various plots and algorithms available for use in estimating future criminal activity.

V. RESULTS

Score Table:

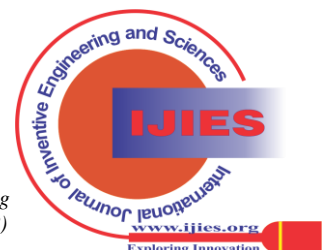
Algorithm	Score
Linear Regression	91%
Logistic Regression	87%
K-Nearest Neighbor	80%
Random Forest Regressor	84%
Random Forest Classifier	91%
Lasso Regression	91%
Decision Tree Classifier	90%
Lasso-KNN	85%
Support Vector Regressor	90%
Lasso-SVR	91%

Given that the output of Lasso-SVR is superior to that of competing algorithms, we adopt a Lasso-KNN strategy to boost the precision of KNN.

Dataset:

The dataset was obtained from the government dataset website by month, but it must be organized and preprocessed prior to usage.

It has columns such as "Month of the previous year," "Major Head," "Minor Head," "End of month," "prior month," "current month," and "Name of month." We heavily rely on the previous year's month, the previous month, the end of the month, the major head, and the name of the month to forecast the value of the current month.



	Act	Major head	Minor head	End Month	Month Previous	Previous month	Current month	Month
0	A - IPC Crime	Murder (Sec.302/303 IPC)	For gain	6.0	4.0	3.0	6.0	1
1	A - IPC Crime	Murder (Sec.302/303 IPC)	Over Property Dispute	4.0	2.0	5.0	4.0	1
2	A - IPC Crime	Murder (Sec.302/303 IPC)	Due to Personal Vendetta or enmity	2.0	1.0	2.0	2.0	1
3	A - IPC Crime	Murder (Sec.302/303 IPC)	Due to Sexual jealousy	3.0	2.0	2.0	3.0	1
4	A - IPC Crime	Murder (Sec.302/303 IPC)	For dowry by burning	0.0	1.0	0.0	0.0	1
...
5468	D.CRIME AGAINST CHILDREN	Probation of Offenders Act	NaN	3.0	0.0	0.0	0.0	9
5469	E. CRIME AGAINST SCHEDULED CASTES /TRIBES BY ...	Murder	NaN	69.0	7.0	3.0	6.0	9
5470	E. CRIME AGAINST SCHEDULED CASTES /TRIBES BY ...	Rape	NaN	161.0	15.0	25.0	13.0	9
5471	E. CRIME AGAINST SCHEDULED CASTES /TRIBES BY ...	Kidnapping	NaN	16.0	0.0	6.0	2.0	9
5472	E. CRIME AGAINST SCHEDULED CASTES /TRIBES BY ...	Offences under the Protection of Civil Rights ...	NaN	4.0	0.0	0.0	2.0	9

5473 rows × 8 columns

Fig 5.1: Crime dataset

KNN Score:

Here, we utilized the KNN method to determine the predicted output and accuracy. Here, we first train our model using the fit method in Python, and then we predict the test data, which accounts for 30 percent of our dataset. 70% of our dataset is used for training purposes. Here, Knn achieves an accuracy of approximately 80%. Here, we must increase the precision of Knn by combining the Lasso and Knn approaches.

```
In [117]: from sklearn.neighbors import KNeighborsRegressor
In [118]: kn=KNeighborsRegressor(8)
In [119]: kn.fit(x_train,y_train)
Out[119]: KNeighborsRegressor(n_neighbors=8)
In [120]: l=kn.predict(x_test)
In [123]: kn.score(x_test,y_test)
0.803787230493467
```

Fig 5.2: Model Score after applying KNN approach

Lasso-KNN Score:

Here, we first use the Lasso method for training, which predicts all training output and stores all predicted results in the g_train variable. Then, we use the KNN approach for re-training, which includes x_train and g_train (predicted Lasso results). Then, we construct a model that accurately predicts the x-test results. This strategy increases the KNN accuracy by 5%.

```
In [38]: from sklearn.neighbors import KNeighborsRegressor
from sklearn.linear_model import Ridge,Lasso
b=Lasso(alpha=2.0)
b.fit(x_train,y_train)
g_train=b.predict(x_train)
In [39]: kn=KNeighborsRegressor(8)
In [40]: kn.fit(x_train,g_train)
Out[40]: KNeighborsRegressor(n_neighbors=8)
In [41]: kn.predict(x_test)
Out[41]: array([-0.09811477, 0.24767629, 0.05188954, ..., 6.90322526,
0.48076657, 0.52521337])
In [43]: kn.score(x_test,y_test)
Out[43]: 0.8461155447033872
```

Fig 5.3: Model Score after applying Lasso-KNN approach

After employing the lasso-KNN method, we were able to successfully improve the KNN accuracy.

Lasso-SVR Score:

Here, we begin by training with the Lasso method, which predicts all training output and stores all predicted results in the g_train variable. Then, we retrain using the SVR method, which includes x_train and g_train (predicted Lasso results). Then, we construct a model that predicts the x-test results accurately. This strategy improves the accuracy of the Lasso by 6%.

```
In [26]: from sklearn.svm import SVR
In [27]: a=SVR(kernel='linear')
In [29]: from sklearn.linear_model import Ridge,Lasso
b=Lasso(alpha=2.0)
b.fit(x_train,y_train)
g_train=b.predict(x_train)
In [30]: a.fit(x_train,g_train)
Out[30]: SVR(kernel='linear')
In [31]: a.predict(x_test)
Out[31]: array([-0.04906938, 0.28591292, 0.18200971, ..., 9.8274206 ,
0.56392885, 0.52531227])
In [33]: a.score(x_test,y_test)
Out[33]: 0.9141598344895
```

Fig 5.4: Model Score after applying Lasso-SVR approach

Predicted output:

The number of victims during the month of August is anticipated here. We are utilizing appropriate supervised learning techniques to make accurate forecasts. When month is set to 8, major head is set to 14, previous year month is set to 5, prior year month is set to 9, and current year up to the end of the month under review is set to 2, the output is 3.47 for Lasso-KNN and 8.22 for SVR. Our proposed algorithm delivers a better result than SVR. As Our proposed algorithm, Lasso-KNN, yields an output of 3.47, which is somewhat bigger than 2 for the current year up to the end of the month under review, but on the other hand, it yields 8.22, which is significantly greater than 2,



therefore, there is a high probability of an overfit error occurring. But we need more accurate results than our proposed method Lasso-KNN can give, so we are using our proposed method Lasso-SVR instead. Our proposed method, Lasso-SVR, gives a result of 6.57, which is better than all of the other algorithms' results. We are still using the two proposed methods, but Lasso-SVR is giving better results than Lasso-KNN. However, by using the Lasso-KNN approach, we were able to improve the accuracy of KNN from 80% to 85%. We use Python, a widely used programming language for statistical analysis, to make predictions about the frequency with which various types of crime will occur in this project.

```
In [213]: a.predict([[8,14,5,9,2]])
Out[213]: array([8.22962592])
```

Fig 5.5: Predicted value of SVR algorithm

```
In [62]: kn.predict([[8,14,5,9,2]])
Out[62]: array([3.4742516])
```

Fig 5.6: Predicted value of Lasso-KNN approach

```
In [168]: kn.predict([[8,14,5,6,2]])
Out[168]: array([[0.5]])
```

Fig 5.7: Predicted value of KNN algorithm

```
In [259]: a.predict([[8,14,5,9,2]])
Out[259]: array([6.57459174])
```

Fig 5.8: Predicted value of Lasso-SVR algorithm

VI. DISCUSSION

This paper talks about crime prediction systems that use machine learning techniques and the ways in which those systems do their job. In machine learning, algorithms like K-Means, SVR, Lasso-KNN, Linear and Logistic Regression, Random Forest, Decision tree, etc. are used to predict events. From these used data analysis algorithms, "K-Means" gives a lower score, so we need to combine it with "Lasso" to raise the score and prevent "overfitting," which happened with the "SVR" algorithm. This happened because the criminal data was reliable and relevant. But K-Means won't work well with noisy data, so we need to preprocess it to deal with missing and NAN values. The algorithm also won't give a better average when there are more clusters. Because of this, K-Means works best when the data is not too noisy and there aren't too many clusters. But when we put lasso and KNN together, we got a prediction that was accurate 85% of the time. On the other hand, we need more accurate results, so we need to use the Lasso-SVR approach to get up to the most accurate output of 6.57, which is much better than the Lasso-KNN approach. In this project, we are using two proposed methods: first, we are increasing the score of KNN, which is up to 85%, and then we are using the best approach to get up to the score of 91%. These are two

possible ways to do something. The area being talked about has grown because crime pattern detection systems now also use image processing.

VII. CONCLUSION

With the use of machine learning technologies, it has become easier to identify relationships and patterns within diverse data sets. This study focuses mostly on predicting the type and number of crimes that may occur. Using the notion of machine learning, we constructed a model with training data sets that underwent data cleansing and modification. When we use Lasso-KNN, the model can predict the type of crime with 0.846% accuracy. When we use Lasso-SVR, it can do the same thing with 0.914% accuracy. Data visualization facilitates data set analysis. The graphs consist of bar, line, and scatter graphs, each with their own distinct qualities. We developed numerous graphs and discovered intriguing statistics that helped us comprehend Indian crime datasets that can aid in identifying elements that contribute to a safer society. Our program provides a framework for viewing the number of crimes and analyzing them using a variety of machine learning algorithms. Using a variety of interactive visualizations, the initiative assists criminal analysts in analyzing these crimes. The interactive and visual feature applications will aid in the reporting and identification of criminal patterns. Clearly, law enforcement agencies may benefit greatly from utilizing machine learning algorithms to combat crime and save mankind.

ACKNOWLEDGEMENT

This paper's author is grateful to Mr. Somen Nayak and Dr. Yogesh Kumar Jakhar, Assistant Professor, for their assistance and support during this independent research endeavor.

DECLARATION

Funding/ Grants/ Financial Support	No, I did not receive.
Conflicts of Interest/ Competing Interests	No conflicts of interest to the best of our knowledge.
Ethical Approval and Consent to Participate	No, the article does not require ethical approval and consent to participate with evidence.
Availability of Data and Material/ Data Access Statement	The data has been extracted from the website of the Indian government (https://data.gov.in/catalog/crime-review-report-2022) and is being pre-processed to improve its format.
Authors Contributions	I am only the solo author of the article.

REFERENCES

1. Lakshman Narayana Vejendla and A Peda Gopi, (2019), "Avoiding Interoperability and Delay in Healthcare Monitoring System Using Block Chain Technology", *Revued' Intelligence Artificielle* , Vol. 33, No. 1, [CrossRef]
2. Gopi, A.P., Jyothi, R.N.S., Narayana, V.L. et al. (2020), "Classification of tweets data based on polarity using improved RBF kernel of SVM" . *Int. j. inf. tecnol.* (2020). [CrossRef]
3. A Peda Gopi and Lakshman Narayana Vejendla, (2019), "Certified Node Frequency in Social Network Using Parallel Diffusion Methods", *Ingénierie des Systèmes d' Information*, Vol. 24, No. 1, 2019, pp.113-117. [CrossRef]
4. Lakshman Narayana Vejendla and Bharathi C R ,(2018), "Multi-mode Routing Algorithm with Cryptographic Techniques and Reduction of Packet Drop using 2ACK scheme in MANETs", *Smart Intelligent Computing and Applications*, Vo1.1, pp.649-658. [CrossRef]
5. Lakshman Narayana Vejendla and Bharathi C R, (2018), "Effective multi-mode routing mechanism with master-slave technique and reduction of packet droppings using 2-ACK scheme in MANETs", *Modelling, Measurement and Control A*, Vol.91, Issue.2, pp.73-76. [CrossRef]
6. Lakshman Narayana Vejendla , A Peda Gopi and N.Ashok Kumar,(2018), "Different techniques for hiding the text information using text steganography techniques: A survey", *Ingénierie des Systèmes d'Information*, Vol.23, Issue.6, pp.115-125. [CrossRef]
7. Crime review report of India ,Open government data platform, NIC, Ministry of Electronics & Information Technology, <https://data.gov.in/catalog/crime-review-report-2022>".
8. Patibandla, R.S.M.L., Veeranjanyulu, N. (2018), "Performance Analysis of Partition and Evolutionary Clustering Methods on Various Cluster Validation Criteria", *Arab J Sci Eng* ,Vol.43, pp.4379-4390. [CrossRef]
9. R S M Lakshmi Patibandla, Santhi Sri Kurra and N.Veeranjanyulu, (2015), "A Study on Real-Time Business Intelligence and Big Data", *Information Engineering*, Vol.4, pp.1-6. [CrossRef]
10. K. Santhisri and P.R.S.M. Lakshmi,(2015), "Comparative Study on Various Security Algorithms in Cloud Computing", *Recent Trends in Programming Languages* , Vol.2, No.1, pp.1-6.
11. Patibandla, R. S. M. Lakshmi et al., (2016), "Significance of Embedded Systems to IoT.", *International Journal of Computer Science and Business Informatics*, Vol.16, No.2, pp.15-23.
12. AnveshiniDumala and S. PallamSetty. (2020), "LANMAR routing protocol to support real-time communications in MANETs using Soft computing technique", *3rd International Conference on Data Engineering and Communication Technology (ICDECT-2019)*, Springer, Vol. 1079, pp. 231-243. [CrossRef]
13. AnveshiniDumala and S. PallamSetty.(2019), "Investigating the Impact of Network Size on LANMAR Routing Protocol in a Multi-Hop Ad hoc Network", *i-manager's Journal on Wireless Communication Networks (JWCN)*, Volume 7, No. 4, pp.19-26. [CrossRef]
14. Paul E. Barrett, J.Hunter, J.T. Miller, J.-C.Hsu, "matplotlib, A portable python plotting package", *Conference, Astronomical data analysis software and systems XIV* , Volume:347.

AUTHORS PROFILE



Prof. Suman Acharya is an assistant professor at the University of Engineering and Management in Jaipur, Rajasthan. He received his BCA in Computer Applications from the Institute of Engineering and Management, Kolkata, West Bengal, and his MCA in Computer Applications from the University of Engineering, Kolkata, West Bengal. He also participated in a number of workshops and faculty development programs. His research focuses on machine learning, data mining, and decision learning, all of which combine to form the discipline of artificial intelligence. Currently, he must shift his research focus to quantum computing in order to apply his expertise in artificial intelligence to a quantum artificial intelligence.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of the Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP)/ journal and/or the editor(s). The Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP) and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.