

Crime Analysis using Supervised Learning

by . Sabyasachi Chakraborty

Submission date: 19-Mar-2019 11:57AM (UTC+0530)

Submission ID: 1095862216

File name: Final_Updated_paper_set_of_crime_analysis.docx (169.88K)

Word count: 3399

Character count: 19407

Crime Analysis using Supervised Learning

Sabyasachi Chakraborty

14th B.Tech. in Computer Application
Vellore Institute of Technology
University in Vellore
Tamil Nadu, India.

Pratik Kumar Dutta

13th B.Tech. in Computer Application
Vellore Institute of Technology
University in Vellore
Tamil Nadu, India

Sahil Gandhe

Master of Computer Application
Vellore Institute of Technology
University in Vellore
Tamil Nadu, India.

Under the guidance of

10th Prof. Jayakumar S
(Assistant Professor)
School of Information Technology and
Engineering
Vellore Institute of Technology
University in Vellore
Tamil Nadu, India.

Abstract

This paper describes that Crime analysis is performed using an ensemble of Naive Bayesian Classification, Random Forest and Support Vector machine algorithms. In this technique, team has calculated and given the result of the accuracy for the provided terms and data. System is able to get prediction automatically with this approach and approximate outcome for a big division of Crime analyzed data and expressions by comparing and deploying efficient algorithms that are significantly better and thereby achieving results in minimal span of time.

Keywords- Crime Analysis and supervised learning, machine learning, Algorithms ensemble and compare of dataset.

I. INTRODUCTION:

Crimes are the increasing danger to the mankind. There are numerous crimes that happens normal interim of time. Maybe it is expanding and spreading at a quick and huge rate. Crimes occur from little town, town to enormous urban areas. Violations are of various kinds – burglary, murder, assault, ambush, battery, false detainment, hijacking, manslaughter. Since Crimes are expanding there is a need to illuminate the cases in a lot quicker way. The Crime exercises have been expanded at a quicker rate and it is the obligation of police division to control and decrease the crimes exercises. Crimes forecast also, criminal recognizable proofs are the serious issues to the police office as there are colossal measures of Crimes occurrence information that exist. There is a need of innovation through which the case explaining could be quicker.

This paper represents about the field of law requirement has developed into a very unpredictable profession with endless zones of claim to fame and ability. Crime investigation could be considered the most current expansion to the field. While the genuine term has been being used since the nineteenth century, crimes examination has a late out come to be viewed as it is very own control in the law implementation field. In any case, since it is so new, there

are as yet numerous issues to be settled encompassing it's across the board usage and use.

Procedure and techniques of crime analysis is the process of recognizing and sorting conclusions communicated in a piece of content, particularly with the end goal to decide if the item is positive, negative or unbiased along with comparing to get better accuracy after the pre processing. Different application like emotion detection and Crime identification, Age prediction, Selection and Multilayer Perception can be done using the techniques 1. Classifier like Naive Bayesian, 2. Random Forest and 3. Support vector machine.

Issue made us to go for an examination about by what method can be workable a crime case made less demanding. Through numerous documentation and cases, it turned out that AI what's more; information science can make the work simpler and quicker. The point of this venture is to make crimes forecast utilizing the highlights present in the dataset. The dataset is removed from the official locales. With the assistance of AI calculation, utilizing python as center we can anticipate the kind of crimes which will happen in a specific region.

Information of particular crime data and its recovery systems mostly go around handling, comparing and sequentially arranging also examining and executing down the truthful information present in it. Details have a target segment in any case, there are some other data encrypted substance which express abstract attributes. These substances are for the most part conclusions, estimations, examinations, dispositions, and crime statistics which shape the center of Crime Analysis (CA). It offers many testing chances to grow new applications, for the most part because of the tremendous development of accessible data on online sources like web journals and interpersonal organizations, Internet, news and media services and sites. For instance, proposals of components proposed by a suggested system can be anticipated by considering contemplation, for an

example, positive or negative and neutral attributes about Dataset by making utilization of Crime Analysis.

Crime data study and process analysis

The arrangement while performing Crime analysis study and process investigation on Crime data set is essentially to order the serially arrange categorized classes precisely and independently. In this ground of research which developed a methodology can propose techniques to prepare a model and after that test it to check its proficiency. As mentioned, main test is to executing crime analysis study in various places of country.

- (1) **Restricted crime data amount:** Restricted amount of data are defined and with only 150 characters close by, minimized articulations are produced, which results sparse arrangement of highlights.
- (2) **Dataset type:** This part defines that it can permit the utilization of table labels and names of state where the crime Occurring, client reference and URLs. These require unexpected preparing in comparison to different data in all set.
- (3) **Consumer range:** Possible of Clients express their idea in an assortment of ways, some utilizing diverse data of criminal Activity in the middle, while others utilizing rehashed words or images to pass on a crime data predictions. Each one of these issues is required to be looked in the pre-processing area. Apart from these, we face problems in feature extraction with fewer features in hand and reducing the dimensionality of features.

II. LITERATURE REVIEW

- a) The authors[1] Agrawal, R., Imielinski, T. and Swami A, in the year of 1993 describes spatial grouping is a movement methodology assembling a great deal for geo referred point-data, $P = \{p_1, p_2, \dots, p_n\}$ in examination area S, into tinier homogeneous parts as a result of continuity (closeness into special extraction). This distinguishes examples with space fixations in expansive space DB. People do qualities, shortcomings. Distinguished space integrated conglomerations were characteristic with fascinating regions (worldwide problem areas and limited abundances) that require further investigation to discover causal variables or conceivable relationships. ARM algorithm for figure full affiliation system fulfilling client indicated least help and least certainty requirements. An Association rule is an articulation as $X \Rightarrow Y (c\%)$, where X can be forerunner and Y is the subsequent, X and Y were sets cause for things in value-based DBs, $X \cap Y = \emptyset$. This is translated with "c% information can fulfill X likewise fulfill Y".

Support form is intended with $P[X \cap Y]$,

Certainty is a gauge with $P[X \cap Y]/P[X]$.

Help are proportion for exchanges fulfill both attributes in quantity for exchanges attach with DB. That certainty was restrictive likelihood of Y given X. Since clients are keen on extensive help and high certainty (solid guidelines (Koperski and Han, 1995)), two edges (least help and least certainty) are utilized for pruning principles to discover solid association rules.

Advantage: Association-rule mining has been a powerful tool for discovering correlations among massive databases. The association-rule mining can easily compute all association rules satisfying user-specified minimum support and minimum confidence constraints.

Disadvantage: Sometimes there are difficulties with efficiency, effectiveness and degree of autonomy in clustering methods, difficulty of cluster shape extraction and lack of adequate correlation measures.

- b) The Authors [2] Oatley and Ewartt in the year of 2004 considers five topographical layers as portrayed. On the off chance that we pinpoint an area inside S, the area will have five related ascribes comparing to the five layers. 5 properties vertically (alluded as a property 3D shape in this paper). Estimations of characteristics turn out to be valid (1) if the area exists in districts (bunches) of comparing layers, false (0) generally. The vertical-see approach endeavours to find intriguing relationship from the entire arrangement of quality shapes. For example, an affiliation rule "**layer(1) \wedge layer(2) \Rightarrow layer(4) (70%)**" is inferred if 70% of trait blocks fulfilling property estimations of layer 1 and layer 2, likewise have the esteem valid in layer 4.

1 The flat view approach overlays every one of the layers into an objective layer and after that endeavors to discover relationship from the objective layer utilizing convergence (covering) regions. The first and second layers converge while the third layer does not cross with the other two. Along these lines, the relationship between the primary layer and second layer is higher than that of the first and third and that of the second and third.

Advantage: Multivariate systems enable scientists to take a gander at connections between factors in a larger route and to measure the connection between factors.

Points of interest of multivariate examination incorporate a capacity to gather a more reasonable picture than taking a gander at a solitary variable.

Disadvantage: Multivariate procedures are perplexing and include abnormal state science that requires a factual program to dissect the information.

- c) The authors [3]Keyvanpour, Mohammad Reza, Mostafa Javideh, and Mohammad Reza Ebrahimi in the year of 2011 gives information removing vital elements related to crime branch investigator and their account mentioned in normal substance. With using procedure, dataset investigation data could normally gone into a DB, with ruling approval workplaces. It occurs moreover associated SOM clustering procedure under degree of wrongdoing investigation examination finally we will use the gathering marks in order for complete the analysis examination coordinating procedure.

Along these lines, a proposed methodology for crime analysis information bunching is spoken to which uses SOM neural system so as to defeat other grouping strategies downsides. The double encrypting causes the noticeable Euclidian partition measure-which is commonly used for constant sorts of components to be useless. That cause of acting parallel sums similar to persistent amounts could prompt misdirect brings about bunching process. Some other separation capacities ought to be utilized which are explicit for accomplishing the likeness between twofold information objects. These capacities ascertain the difference between two articles as their comparing separation. The separation (disparity) between two articles can be determined by conditions given.

$$D(I, j) = 1 - S(I, j)$$

S and D in that order address components for closeness and uniqueness flanked by 2 twofold progressions given in formula. As demonstrated the resemblance with in abit game plans formula attribute can be assessed by the accompanying conditions: -

Coefficient of normal Matching: $s(i, j) = (m+s) / (m+n+l+o)$

Coefficient by Rao's: $s(i, j) = m / (m+n+l+o)$

Jaccard: $s(i, j) = m / (m+n+l)$

Advantage System for crime analysis examination helps government organization by burrowing and dissecting records of the budgetary exchange to assemble designs that can distinguish tax evasion or criminal exercises.

Disadvantage: Complex and include abnormal state science that requires a measurable program to investigate the information.

- d) The Authors [4]Al-Janabi, Kadhim B. Swadi in the year of 2011 considers cluster analysis is an essential human action which revels from youth when figure out how to recognize creatures and plants, and so forth by

constantly improving intuitive bunching plans. It is generally utilized in various applications including design acknowledgment, information examination, picture preparing, and statistical surveying and so on. Bunch precision can be improved to catch the neighborhood connection structure by partner each group with the mix of the measurements as autonomous weighting vector and subspace range which is installed on it. Ongoing advancements in crime analysis control applications go for embracing information mining systems to help the procedure of crime analysis examination. Late research takes a shot at crime analysis examination Includes-Adderley and Musgrove put forth a concentrated effort Arranging Map (SOM) to associate liable gatherings for authentic rapes. Starting late, Ozgulproposed a novel desire exhibit Crime Prediction Model to anticipate positions of not fixed dread monger events in characteristics of wrongdoing investigation information which can be in zone, modusoperandi qualities.

$$Entropy(S) = p_+(-\log_2 p_+) + p_-(-\log_2 p_-) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

P=probability of positive occurrence

Advantage: The pre prepared information were utilized to discover diverse wrongdoing and criminal patterns and practices, and violations and offenders were assembled into bunches as indicated by their vital qualities.

Disadvantage: Complex and include abnormal state arithmetic that requires a measurable program to dissect the information.

- e) The Authors [5]ShyamVaran Nath in the year of 2006in this paper presents theutilization of bunching calculation for an information mining way to deal with assistance distinguish the crime analysis examples and accelerate the way toward tackling crime analysis. We will see k-implies bunching with certain improvements to help during the time spent distinguishing proof of crime analysis designs. We connected these methods to genuine crime analysis information from a sheriff's office and approved our outcomes. We likewise use semi-administered learning method here for information revelation from the crime analysis records and to help increment the prescient precision.

Bunch (of crime analysis) has an extraordinary significance and alludes to a geological gathering of crime analysis, for example a ton of crime analysis in a given topographical area. Such bunches can be outwardly spoken to utilizing a geo-spatial plot of the crime analysis overlaid on the guide of the police locale. The thickly populated gathering of crime analysis is utilized to outwardly find the 'problem areas' of crime analysis. In any case, when we discuss bunching from an information mining point of view, we allude to comparable sorts of crime activity in the given geology of intrigue. Such groups are helpful

in recognizing a crime analysis design or a crime analysis binge. Some outstanding instances of crime analysis designs are the DC expert rifleman, a sequential attacker or a sequential executioner. These crime analysis may include single suspect or might be submitted by a gathering of suspects.

Advantage: This simple to actualize information mining system works with the geospatial plot of crime analysis and improves the profitability of the investigators and other law authorization officers. This can likewise be connected for counter fear mongering for country security

Disadvantage: A portion of the impediments of our examination incorporate that crime analysis design investigation can just support the analyst, not supplant them. Information mining is touchy to nature of info information that might be mistaken; have missing data, be information passage blunder inclined and so forth.

- f) The Authors [6] Alves, Luiz GA, Haroldo V. Ribeiro in the year of 2018 uses a random forest regression process is to anticipate criminal activity analysis measure their impact for urbana pointers in manslaughters. Users technique has till 97% of exactness in wrongdoing investigation desire and hugeness with urbana pointers are situated and clustered with social events for proportionate effect by solid beneath possibly transform with information test examined. Outcomes decide position for significance for urbana markers which anticipate criminal activity analysis, uncovering joblessness and ignorance were the most essential factors of depicting murders in Brazil citizens urban areas.
- g) The Authors [7] Shiju Sathyadevan, Surya Gangadharan S in the year of 2014 describes examination of Crime prediction and analysis and aversion is an efficient methodology for recognizing and dissecting examples and patterns in crime prediction analysis. Our framework can anticipate locales which have high likelihood for crime prediction event and can envision crime prediction and analysis inclined zones. With the expanding coming of electronic frameworks, crime analysis information analytics team can help the Law requirement officers to accelerate the way toward tackling violations. Utilizing the idea of information mining we can remove already obscure, helpful data from unstructured information. Here we have a methodology between software engineering and criminal equity to build up an information mining strategy that can help illuminate crime prediction quicker. Rather than concentrating on reasons for criminal active event like criminal foundation of wrongdoer, political hostility and so forth we are concentrating

fundamentally on crime analysis components of every day.

Advantage: This setting is propelled by numerous cases in which there exist laws that refuse a choice that is mostly founded on segregation. Guileless use of AI strategies would result in colossal fines for organizations.

Disadvantage: Cannot deal with expansive and complex dataset precisely.

III. BACKGROUND

A. Naive Bayesian classifier

We can say Naive Bayesian classifiers defines that this classifiers are considering the characterization job from a Statistical perspective. The beginning stage is that the likelihood of a class is given by the back likelihood or probability given a preparation report. Here alludes to the majority of the content in the whole preparing set. Here is given by , where is the aspect (word) of record .

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} \quad [1]$$

Probability = probability of two objects
A = first object
B = Second object

B. Random Forest

The process manufactures numerous choice trees and consolidates them to get a progressively precise and stable forecast

Using Gini Split / Gini Index

$$G = 1 - \sum_{i=1}^c (P_i)^2 \quad [2]$$

P = probability of Positive Occurrence
G = Gini index

C. Support Vector Machine

The process is an administered AI calculation which can be utilized for both grouping and relapse difficulties

$$h_{w,b}(x) = g(w^T x + b) \quad [3]$$

w, b = two attributes

Can be generalized to more than two values

IV. MODULES

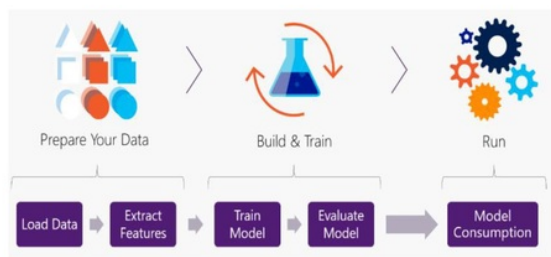
- Data Collection

- Data Pre-processing
- Algorithm extraction:
 - a) Naive Bayesian prediction
 - b) Random Forest prediction
 - c) Support Vector Machine (SVM) prediction
- Features of algorithm extraction
- Compare the result based on accuracy
- Analyze the performance

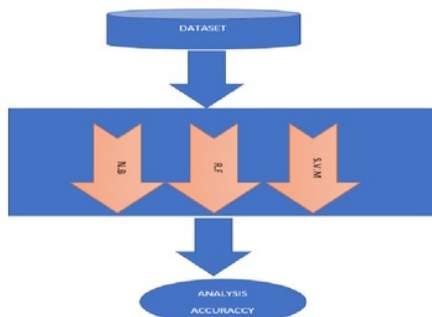
V. PROPOSED MODEL

- The model is an ensemble model of three algorithms.
- In this model we analysis data based on Bayesian classifier, R Forest and (SVM).
- All three classifiers are supervised which requires the class label for training a testing the data.
- The advantage of using this model is to improve the accuracy prediction.
- Weighted majority voting is used to predict the final accuracy from all the classifiers.

A. Architectural diagram of ensemble algorithm –



Simplified diagram of process –



B. Performance Metrics –

- **Accuracy:** Result of calculation or specification conforms to the correct value or a standard of three algorithms. Accuracy is used as a metric for evaluating the effectiveness of a classifier.

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+FN+TN} \text{ ----- [4]}$$

- **Precision:** Refinement in three algorithms for calculation, or specification, especially as represented by the number of digits given. By definition we can elaborate that Precision can estimate the accuracy of one classifier. The higher precision goes, that implies less false positives. On other hand a lower precision implies all the high false positives. This shows regularly inconsistent comparing with recall part, as a simple method to enhance precision is to diminish recall.

$$\text{Precision} = \frac{TP}{TP+FP} \text{ ----- [5]}$$

- **Recall:** Recall section defines that deals with the culmination, or affect ability, of a classifier and measures for output. Higher recall implies less false negative, while lower recall review implies all the more false negatives. Enhancing recall can regularly diminish accuracy since it persuades progressively harder to be exact as the demo space increments.

$$\text{Recall} = \frac{TP}{TP+FN} \text{ ----- [6]}$$

VI. RESULTS

Followings are the individual outputs of three algorithms by comparing all and accuracy of each:

- *Naive Bayesian Classifier –*

```
Accuracy for gaussian Naive Bayes : 0.761608040201005
Recall for gaussian: 0.692
Precision for gaussian: 0.9117998148278733
```

- *Random Forest –*

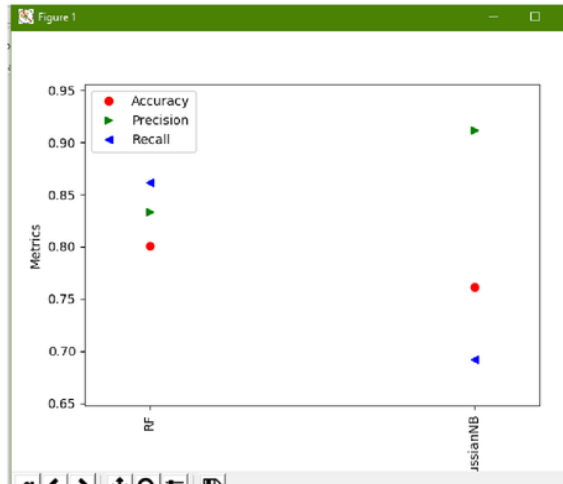
```
Accuracy for RandomForestClassifier is:- 0.8012537688442212
Precision for RandomForestClassifier is 0.8333878362200029
Recall for RandomForestClassifier is 0.8615999999999999
```

- *Support Vector Machine –*

```
Accuracy for SVM is: 0.8072713567839195
Precision for SVM is: 0.8549456169037868
Recall for SVM is 0.8448
```

- *Comparing algorithms –*

Comparing two, Bayesian classifier and R F in graph format:



VII. CONCLUSION

In this paper, Crime Analysis is performed by collecting the Crime data and individual pacifiers such as three algorithms mentioned in the project (SVM, RF, and Bayesian).

The individual classifier accuracy of 80% is obtain for two classifiers which are SV Machine and R Forest algorithms and 76% of Bayesian Classifier however in the future, an ensemble of three algorithms is deployed for better accuracy.

VIII. REFERENCES

- [1] R. Agrawal, T. Imieliński, and A. Swami, "Mining association rules between sets of items in large databases," in *Acm sigmod record*, 1993, vol. 22, no. 2, pp. 207–216.
- [2] G. C. Oatley, J. Zeleznikow, and B. W. Ewart, "Matching and predicting crimes," in *Applications and Innovations in Intelligent Systems XII*, Springer, 2005, pp. 19–32.
- [3] M. R. Keyvanpour, M. Javideh, and M. R. Ebrahimi, "Detecting and investigating crime by means of data mining: a general crime matching framework," *Procedia Comput. Sci.*, vol. 3, pp. 880–880, 2011.
- [4] K. B. S. Al-Janabi, "A proposed framework for analyzing crime data set using decision tree and simple k-means mining algorithms," *J. Kufa Math. Comput.*, vol. 1, no. 3, pp. 8–24, 2011.
- [5] S. V. Nath, "Crime pattern detection using data mining," in *2006 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology Workshops*, 2006, pp. 41–44.
- [6] L. G. A. Alves, H. V. Ribeiro, and F. A. Rodrigues, "Crime prediction through urban metrics and statistical learning," *Phys. A Stat. Mech. its Appl.*, vol. 505, pp. 435–443, 2018.
- [7] S. Sathyadevan, S. Gangadharan, and others, "Crime analysis and prediction using data mining," in *2014 First International Conference on Networks & Soft Computing (ICNSC2014)*, 2014, pp. 406–412.

Crime Analysis using Supervised Learning

ORIGINALITY REPORT

8%

SIMILARITY INDEX

8%

INTERNET SOURCES

7%

PUBLICATIONS

%

STUDENT PAPERS

PRIMARY SOURCES

1

www.geocomputation.org

Internet Source

2%

2

www.grdjournals.com

Internet Source

1%

3

sci2s.ugr.es

Internet Source

1%

4

export.arxiv.org

Internet Source

1%

5

www.sersc.org

Internet Source

1%

6

core.ac.uk

Internet Source

<1%

7

Mehedee Hassan, Mohammad Zahidur Rahman. "Crime news analysis: Location and story detection", 2017 20th International Conference of Computer and Information Technology (ICCIT), 2017

Publication

<1%

Irina Matijosaitiene, Peng Zhao, Sylvain

8

Jaume, Joseph Gilkey Jr. "Prediction of Hourly Effect of Land Use on Crime", ISPRS International Journal of Geo-Information, 2018
Publication

<1 %

9

Malith Munasinghe, Harsha Perera, Shanika Udeshini, Ruwan Weerasinghe. "Machine Learning based criminal short listing using Modus Operandi features", 2015 Fifteenth International Conference on Advances in ICT for Emerging Regions (ICTer), 2015
Publication

<1 %

10

thesai.org
Internet Source

<1 %

11

Mehmet Sait Vural, Mustafa Gok, Zeki Yetgin. "Generating incident-level artificial data using GIS based crime simulation", 2013 International Conference on Electronics, Computer and Computation (ICECCO), 2013
Publication

<1 %

12

link.springer.com
Internet Source

<1 %

13

journals.plos.org
Internet Source

<1 %

14

www.nature.com
Internet Source

<1 %

Exclude quotes On

Exclude bibliography On

Exclude matches

< 10 words