# Crime Pattern Detection, Analysis & Prediction

Sunil Yadav, Meet Timbadia, Ajit Yadav, Rohit Vishwakarma and Nikhilesh Yadav
Department of Information Technology,
University of Mumbai, Shree L.R Tiwari College of Engineering, Thane, India
sunvan.15@gmail.com, meettimbadia@gmail.com, ajitdmyadav@gmail.com,
rohitvish29@gmail.com,nikhileshyadav29@gmail.com

**Abstract: Crimes are a social irritation and cost our society deeply in several ways. Any research that can help in solving crimes quickly will pay for itself. About 10% of the criminals commit about 50% of the crimes [9]. The system is trained by feeding previous years record of crimes taken from legitimate online portal of India listing various crimes such as murder, kidnapping and abduction, dacoits, robbery, burglary, rape and other such crimes. As per data of Indian statistics, which gives data of various crime of past 14 years (2001-2014) a regression model is created and the crime rate for the following years in various states can be predicted [8]. We have used supervised, semi-supervised and unsupervised learning technique [4] on the crime records for knowledge discovery and to help in increasing the predictive accuracy of the crime. This work will be helpful to the local police stations in crime suppression.**

*Keywords: K-means, Naive Bayes, Crime prediction, Regression, Crime suppression, Apriori, Association Rule*

## I. INTRODUCTION

The crime rates accelerate continuously and the crime patterns are constantly changing [2]. As a result, the behaviors in crime pattern are difficult to explain. This paper illustrates how social development may lead to crime prevention. The aim is to provide a comprehensive review of theory and research with respect to the prevention of the crime in the society and to implement different data analysis algorithms which address the connections between crime and its pattern. The data for the project are collected from the legitimate government sources [7]. The data was converted to .csv format upon which pre-processing of the data was performed. Technologies used for mining various crime pattern and analysis are Weka Tool and R Tool.

*Weka Tool:*

Weka is a collection of machine learning algorithms for data mining tasks. The algorithms can either be applied directly to a dataset or called from your own Java code. Weka contains tools for data pre-processing, classification, regression, clustering, association rules, and visualization. [11].

*R Tool:*

R is a language and environment for statistical computing and graphics. It is a GNU project which is similar to the S language R provides a wide variety of statistical (linear and nonlinear modeling, classical statistical tests, time-series analysis, classification, clustering and graphical techniques, and is highly extensible [12].
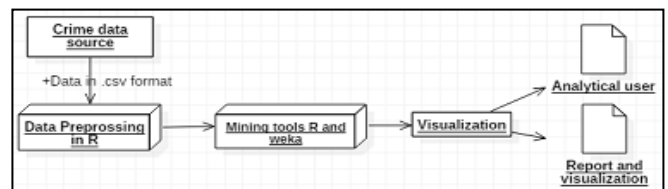


Fig 1.Block diagram(process of mining)

The four algorithms which were executed on the dataset with the help of Weka Tool and R Tool are:

i. *Association Mining (Aprori)*
ii. *Clustering(k-Means)*
iii. *Classification Techniques (Naive Bayes)*
iv. *Correlation & Regression.*

## II. ARCHITECTURE AND WORKING

The dataset embraces number of people arrested and number of crimes committed along with various other attributes. Here we are primarily using four data mining algorithms for analysis of crime and to find hidden patterns of crime in India.



Fig 2.Data set for K-mean

The data mining techniques used are as follows:-

1: *Association mining (Apriori Algorithm) along with*

*Clustering (k-mean):-*

The paper tends to help specialist in discovering patterns, trends, making forecasts, finding relationships and possible explanations, mapping criminal networks and identifying possible suspects [3].

K-means is used to create number of clusters according to values high and low.

Two clusters are to be created:

    Cluster 1: High number of people involved in crime.

    Cluster 2: Low number of people involved in crime.

The data is first imported from the file named as book.csv into weka for preprocessing, later k-means is applied on this data set using the same graphic user interface (GUI) of Weka.
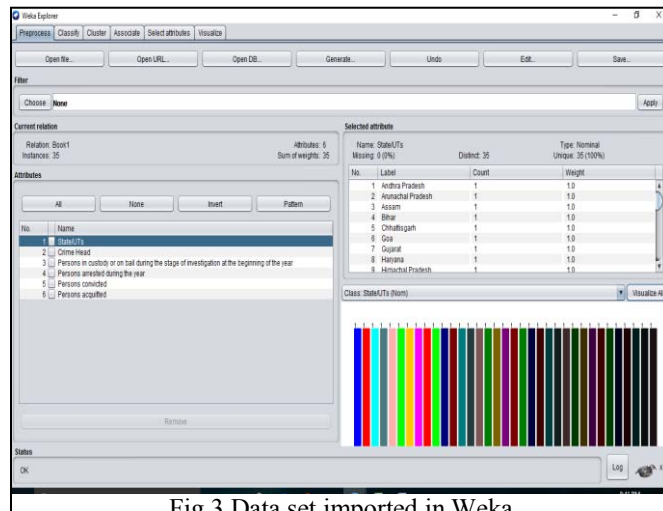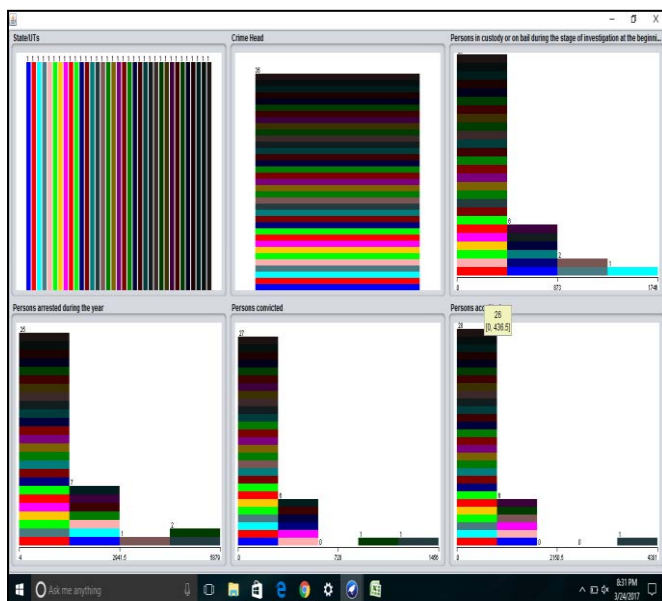


Fig 3.Data set imported in Weka

The imported data set is visualized with class attribute being STATE/UT. The visualization chart shows the distributing of attribute STATE/UT with other five attributes in the data set, each color in the visualization chart represents a particular state



of India.

Fig 4.Visualization of dataset in Weka

The result of the k-means will acts as an input to the Apriori algorithm for discovering the association among a number of



Other attributes.

Fig 5. Result of k-means

*Apriori Algorithm:* - Apriori algorithm [10] is a type of association mining, used to find frequent item set.

The result obtained after k-means is used as dataset for Apriori as the clusters are now divided into high and low value now we



can get the association between various attributes.

Fig 6. Data set for Apriori

The dataset consists of attributes such as person arrested, person convicted, person acquitted. The "H" in the dataset represents the high number of person involved in a particular crime.

As apriori is to be applied on the output of k-means after some data preprocessing, the data set for apriori  is imported into weka and association will be found using weka
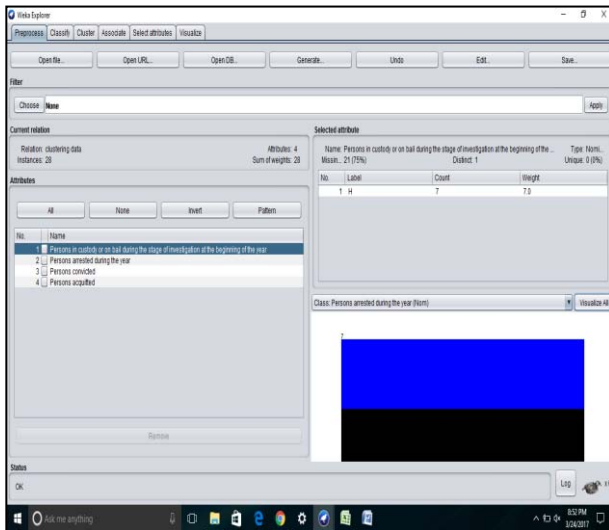
Fig 7. Data set imported in weka

The imported dataset is visualized in weka, the visualization chart shows the distribution of crime as high or low of particular attributes with class attribute which is persons arrested during the year.
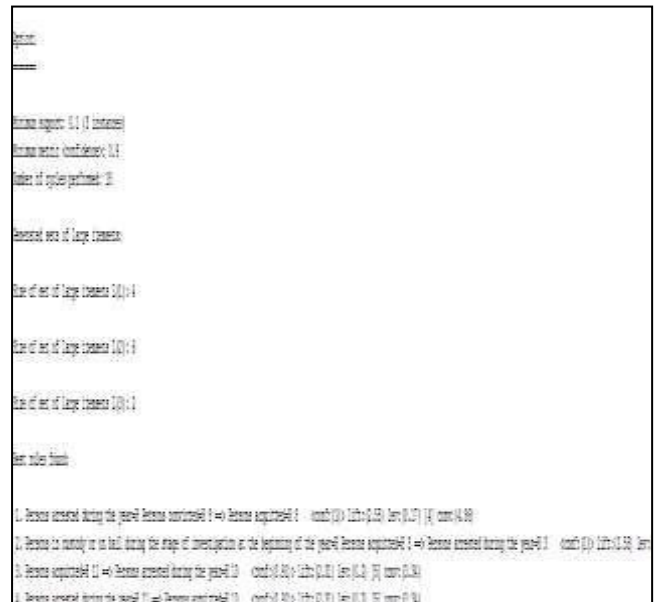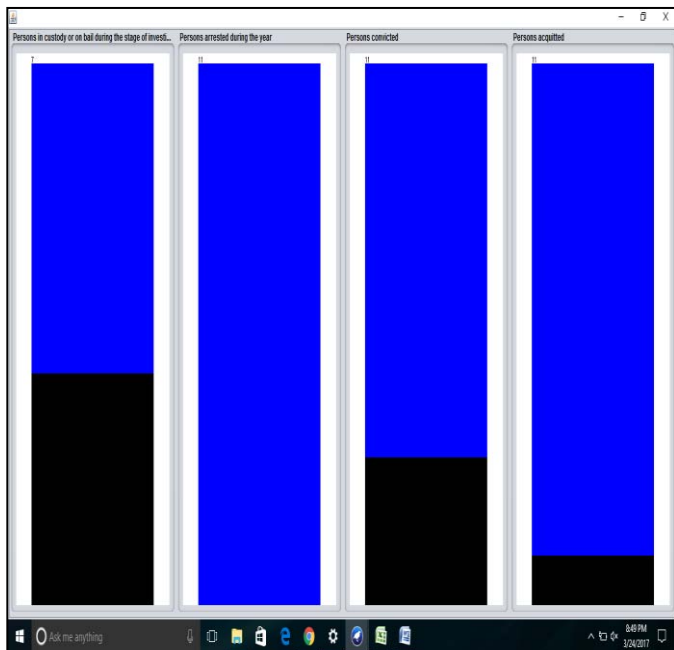

Fig 8. Visualization of Apriori dataset

The blue region in the chart represents the high crime and black region represents the low crime of particular attribute in the dataset.

The result of Apriori showed an association between the person arrested during the year and person acquitted in the same year. This result states that if person arrested are more, then person released are also more, hence more people turn out to be innocent..


Fig 9. Result of Apriori

*2. Naive Bayes Algorithm:-*

Classification is one of the classic data mining techniques, which is used to classify each item in a set of data into one of Predefined set of classes or groups [6]. The idea is to define the Criteria use for the segmentation of the whole database, once this is done, individual dataset can then fall into one or more groups naturally.

With the help of classification, existing dataset can easily be understood and it also helps to predict how new individual dataset will behave based on the classification criteria. Data mining creates classification models by observing already classified data and finding a predictive pattern among those data. [1]


Fig 10. Data set for Naive Bayes

Naive Bayes is a type of classification algorithm used for prediction it works on Bayesian principle, here we have provided a data set for prediction of the number of crime

committed by a particular age group. The sample dataset includes:
STATE/UT-
*AP: Andhra Pradesh*
*ARP: Arunachal Pradesh*

AGE GROUP-
*s-tw: 7-12*
*tw-s:12-16*
NO-CRIMES-
*z-th:0-3*
*f-six:4-6*
*s-nine:7-9*

```
=== Confusion Matrix ===
a b c d e <-- classified as
1 0 0 0 0 | a = z-th
0 0 0 0 0 | b = s-nine
0 0 0 0 0 | c = f-six
0 0 0 0 0 | d = sixteen-eighteen
0 0 0 0 0 | e = thirteen-fifteen
```

Test case for Data set:-

Test case used is of the state Andhra Pradesh where crime type is murder, gender is male, and age group is between 7-8 years

Then number of crime committed by that age group is predicted as follows:

@relation Naive-test

@attribute STATE/UT {AP, ARP} @attribute CRIME {Murder, Rape} @attribute GENDER {m, f} @attribute AGE {s-tw, tw-si}

@attribute NO-CRIMES {z-th, s-nine, f-six, sixteen-eighteen, thirteen-fifteen}

@data

AP, Murder, m, s-tw, z-th

The most probabilistic answer for test case having:

State: Andhra Pradesh
Crime: Murder
Gender: Male
Age Group: 7-8 years

Result:
Number of Murder committed: 0-3

*3. Correlation & Regression:-*

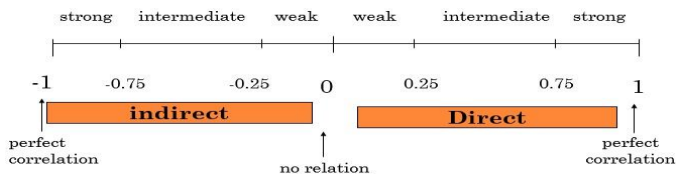Correlation is a statistical technique used to determine the degree to which two variables are related.


Fig 11: Correlation

If the result of correlation is 1 that means there is perfect relation between the two attributes, if result is 0 then there is no relation between the two attributes, hence there must be strong relation between the attributes to get significant result.

The aim of linear regression is to model a continuous variable $Y$ as a mathematical function of one or more $X$ variable(s), so that we can use this regression model to predict the $Y$ when only the $X$ is known. This mathematical equation can be generalized as follows:

$$Y = \beta_1 + \beta_2 X + \epsilon \qquad (1)$$

Where, $\beta_1$ is the intercept and $\beta_2$ is the slope. Collectively, they are called *regression coefficients*. $\epsilon$ is the error term; the part of $Y$ the regression model is unable to explain [13].
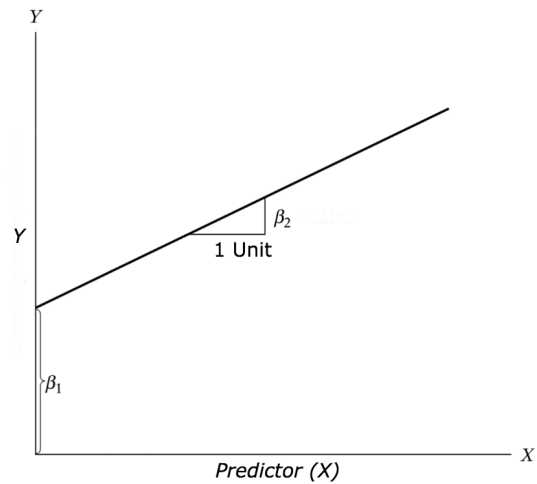

Fig 12: Sample regression diagram

The dataset belong to rape Crime against Women of all the states of India.


Fig 13: Data set for correlation and Regression

The data set of crime rape contains attributes such as number of Trial completed, number of persons convicted during the

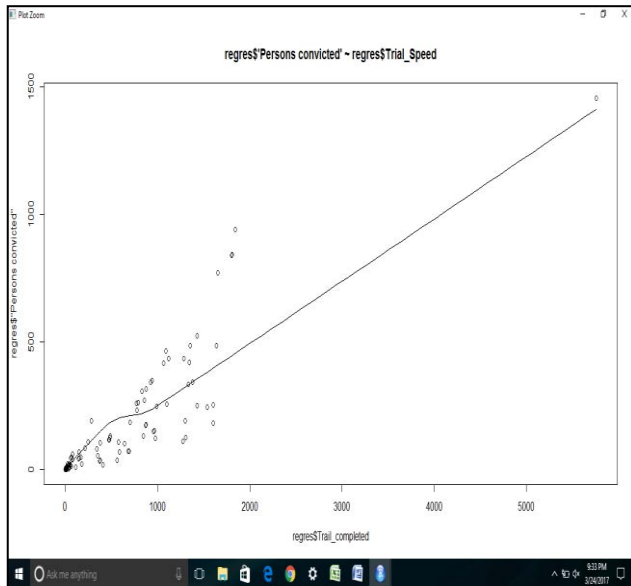trails as well as number of persons acquitted on the charge of crime rape.



Fig 14: Scattered plot of data set

Result of correlation:
X= Trials completed during year. Y=person convicted.

Correlation: 0.98 (Strong relation) There is a strong relation between X & Y.

This means that there is a significant relation between Trials Completed during the year and person convicted on rape charge.



Fig 15: Result of correlation and regression

From output,
$$Y = \beta_1 + \beta_2 X + \epsilon \qquad (2)$$

Here if x is 10 i.e. if trail completed is 10 then

Y (Person convicted) = 0.2767*10-0.3173=2.449     (2.1)

Therefore for every 10 rape case trials completed only 2.5~ 3 people are convicted of the rape charges.

Thus by regression we can predict the number of people who acquitted the crime against the number of Trials completed during the year.

## III.  CONCLUSION AND FUTURE SCOPE

The biggest hurdle in the project was data acquisition and data staging. As a future scope extension of crime detection and analysis will be to generate the crime hot-spots that will help in deployment of police at most likely places of crime for any given window of time, to allow most effective utilization of police resources. The developed model will reduce crimes and will help the crime detection field in many ways that is from arresting the criminals to reducing the crimes by carrying out various necessary measures.
As in near future these methods can be applies on full data set which consists of 42 crime heads having 14 attributes to them, thus when analyzed could provide more unexpected dependencies of attributes over each other.

*Geospatial in Crime Pattern Detection:*

Crime is neither systematic nor entirely random [5].The system could be enriched to support crime mapping over India through which terrestrial model can be created declaring various crimes and the degree of such crimes performed. This terrestrial model would help us to compare the various crime rates in the diverse states of India and to cultivate new strategies to abate the crime rate in that particular area.

## IV.    REFERENCES

[1] Mugdha Sharma, *Z-crime: A data mining tool for the detection ofsuspicious criminal activity based on decision tree*, IEEE, 2014,ISBN:978-1-4799-4674-7/14

[2] Ubon Thansatapornwatana, *A Survey of Data Mining Techniques forAnalyzing Crime Patterns* Second Asian Conference on DefenseTechnology ACDT, IEEE, 2016, ISBN: 978-1-5090-2258-8/16

[3] Dr.Zakaria Suliman Zubi, Ayman Altaher Mahmmud, *Using DataMining Techniques to Analyze Crime patterns in the Libyan National Crime Data*, Recent advances in image, audio and signal processing.ISBN: 978-960-474-350-6

[4]Shiju Sathyadeven, Deven M.S, Surya Gangadharan. S, Crime Analysisand prediction using data mining, IEEE, 2014

[5]Chung- Hsien Yu, Max W. Ward, Melissa Morabito, Wei Ding *Crimeforecasting using data mining technique,*11[th]IEEE

InternationalConference on Data mining workshop, IEEE, 2011,ISBN:978-0-7695-4409-0/11

[6]Fatih OZGUL, Claus ATZENBECK, Ahmet CELIK, Zeki ERDEM, *Incorpating data sources and methologies for crime data mining,* IEEE,2011, ISBN: 978-1-61284-4577-0085-9/11

[7]National Crime Records Bureau Website. [online]

[8]Arunima S. Kumar, Raju K. Gopal, *Data Mining Based CrimeInvestigation System: Taxonomy and Relevance*, IEEE, 2015, ISBN: 978-1-4799-8553-1/15

[9]Shyam Varan Nath, *Crime Pattern Detection Using Data Mining,* International Conference on Web Intelligence and Intelligent Agent (WI-IAT 2006 Workshops), IEEE, 2006, ISBN:0-7695-2749-3/06

[10]Jiawei Han, Micheline Kamber, Jian Pei, *Data Mining: Concepts and Techniques 3$^{rd}$ ed.*

[11]http://www.cs.waikato.ac.nz/ml/weka/.[online]

[12]https://www.r-project.org/about.html.[online]

[13]http://r-statistics.co/Linear-Regression.html.[online]