# Prediction of Crime Occurrence using Multinomial Logistic Regression

**R. Rajadevi, E. M. Roopa Devi, S. Vinoth Kumar**

*Abstract: In order to uncover hidden patterns and correlations, data analysis examines large amounts of data. Analysis of crime isa systematic approach to the identification and analysis of crime patterns and itstrends. This plays a role in the planning of problems with crime and in formulating strategies for crime prevention. Instead of focusing on causes of crime such as criminal offender background, this work focuses primarily crime factors happened on every day. This work can predict the category of crime that has a higher likelihood of occurrence in those areas and can visualize in the form of histogram and heat map by category of crime, crime by day of week and month. The study depends on a lot of variables like class, latitude, longitude, etc. For forecast, the multinomial logistic regression method is used. For weekdays, the district and the hour of the accident are used as predictors.This algorithm is used because its target variable has more than two values and no ordering in the response variable.This provides greater efficiency for handling datasets with multi class labels. This forecast can be helpful in predicting the occurrence of crime in vulnerable areas, which in turn minimizes the crime rate by providing the patrol in those areas.*

*Keywords: Data Analytics, Prediction, Regression ,Machine Learning.*

## I. INTRODUCTION

Big data analytics involves collecting data from different resources to manipulate and then finally deliver as useful products to the organization.

It is useful to integrate raw data acquired from different sources into a data item, forms the core of Big Data Analytics.There are two methodologies in Data Analytics : Exploratory Data Analytics (EDA) and Confirmatory Data Analytics (CDA). EDA is an approach to analyzing data sets to summarize their main characteristics, often with visual methods. We can just see what the data tell us beyond the formal modeling task. In EDA data are explored which can provide information about the numbers of factors required to represent the data. CDA is a multivariate statistical procedure through which we can test how well the measured variables represent the construct.

**Revised Manuscript Received on January 30, 2020.**
∗ Correspondence Author
   **R. Rajadevi∗,** Department of Information Technology, Kongu Engineering College, Perundurai, India. E-mail: rajdevi@kongu.ac.in
   **E. M. Roopa Devi**, Department of Information Technology, Kongu Engineering College, Perundurai, India. E-mail: roopadevi@kongu.ac.in
   **S. Vinoth Kumar**, Department of Information Technology, Kongu Engineering College, Perundurai, India. E-mail: vinoths@kongu.ac.in

CDA is a tool which can be useful to confirm the measurement theory.

There are five Characteristics which is the building blocks of an efficient data analytics solution: Accuracy, Completeness, Consistency, Uniqueness and Timeliness. There is an another characteristic in data analytics called as Data Visualization which describes the presentation of abstract information in graphical form.

It allows users to spot patterns, trends, and correlations that otherwise might go unnoticed in traditional reports, tables, or spreadsheets. There are two basic types of Data Visualization: Exploration and Explanation. By using these categories we have many ways to make data can be visual. The most common types of data visualization are Heat map Cartogram, Choropleth, Dot Distribution Map, Connected Scatter Plot, Polar Area Diagram, Time Series, Pie Chart, Histogram, Scatter Plot, Dendrogram, Ring Chart, Tree Diagram, Alluvial Diagram, Node-Link Diagram, Matrix. A heat map is a two-dimensional representation of data in which values are represented0020by colors. It provides an immediate visual summary of information. It provides easy understanding of complex data sets.

## II. RELATED WORK

. The existing system deals with large set of data and it consist of centralized database. Running Algorithm like Multinomial Logistic Regression has higher time complexity. It consists of 39 categories of crime but classification in somehow difficult. System will consists of poor accuracy and replicated values which will lead to large time consuming. As the system is centralized and it does not distribute the task or data, the retrieval and processing of the system consumes large amount of time. The analysis task is more complex and identifying error rate is difficult.

## III. PROBLEM STATEMENT

The existing method is more complex for analysis and also it has the complex structure which provides complex view to the users. The designers of the System felt so difficult for giving such proper working model. It has large space complexity and time complexity. It seems to be difficult to predict the crime and to process the data from the records of crime. It neither displays the result in the pictorial form nor in the comparative manner. The algorithms such as Random Forest and Naive Bayes has higher complexity in both time as well as space The huge data makes higher complexity in prediction

## IV. PROPOSED SYSTEM

The proposed system is to identify and visualize the occurrence of the crime with higher probability of critical areas. These results are used to predict crime rates in sensitive areas .. The analysis depends on several factors such as latitude and longitude etc.., Data are collected, classified and visualized using graphs. Multinomial logistic regression algorithm is used for prediction. Day of week, District and hour of the incident are used as predictors.

This reason behind in selecting the algorithm is because there are more than two values for target variable and response variables are not ordered. Then the main objective is to predict the area with higher crime rate and to provide security so the crime rate will be reduced. The System has various modules involving Data Visualization and Classification. The data source is the static dataset which is going to be used as a source for the data analysis process. R Studio is used to analyze and visualize the data. The dataset is stored as a csv file for processing which is the first step that has to be done for data processing.
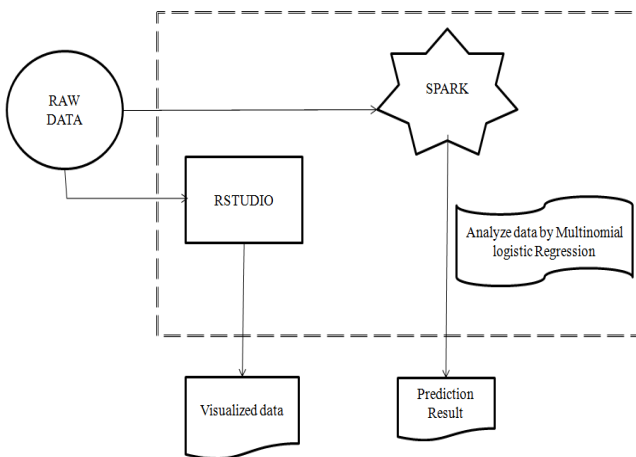


**Fig.1.1 Framework to perform classification**

Import the dataset in the program for training. The category of crime will be removed and the test dataset will be loaded as a csv file. The dataset will be trained in the desired model by multinomial logistic regression. After running the model for several iterations the output file will be generated with the probabilistic values for various crime categories. This system design consists of all the modules which are going to process the datasets by analytics. The raw data which consists of test data and training data. Working directory is set and the data sets are loaded. The data set is given to spark environment and the day of week, district of the crime and the hour of the incident are used as predictors from the dataset. The data is pre-processed by using the methods that deal with the missing values to maximize the crime rate prediction.

The multinomial regression model learns from the parameters weekday, when the crime occurred regularly, city, hour.

### Multinomial Logistic Regression

Relationship of the variable in the dataset is identified by using regression techniques. The categorical variables based on the independent variables can be categorized using Logistic regression .For example to predict whether user will buy a product or not (YES/NO) .Binomial logistic regression model is further modified to categorize the dependent variable which contains nominal values with more than two level using multinomial logistic regression. The dataset contains 39 categorical values like bribery, gambling, runaway, fraud, vehicle theft etc. A multinomial Logistic regression technique is applied to predict the crime based on the input variables.

$$Y1 = Ln(\frac{p2}{p1}) \tag{1}$$

$$Y2 = Ln(\frac{p3}{p1}) \tag{2}$$

$$P1 = 1/(1 + e^y1 + \cdots e^y24) \tag{3}$$

$$P2 = e^y1/(1 + e^y1 + \cdots e^y24) \tag{4}$$

$$Pn = e^y(n-1) /(1 + e^y1 + \cdots e^y2 \tag{5}$$

The above equation No. 3 to 5 is used to find the probabilities The features in the dataset are given as input to model. The values in the features are converted to numerical values. The linear regression model make use of all the features $x1,x2,x3\ldots xn$ and W is initial weight for all the features $W=[w1,w2,w3\ldots.wn]$. The model generated using the features and weights like $w1*x1, w2*x2,w3*x3 \ldots\ldots wn*xn$. During the training phase the initial weight is updated. Softmax function is used to find the probability of each feature and it depends on the updated weights. The one hot encoding is applied for conversion of categorical values to binary format. A matrix is constructed with number of input features and class variable. Distance is calculated using entropy function using Softmax and matrix. During training the expected weights will be obtained .Prediction is done using the weights obtained during training phase. The above procedure is repeated for each record in the training data set until loss function value is minimal. After completing the all the iteration weights obtained it utilized for predicting the expected class. This regression model handles larger data sets with higher efficiency. Evaluation metrics are applied to determine the efficiency of the algorithm.

## V. RESULT AND DISCUSSION

The San Francisco crime data is collected from the Kaggle website. The algorithm chosen for this solution, is a variation of multinomial logistic regression, a classification model based on regression where the dependent variable (what is used as predictor) is categorical (opposite of continuous), implemented in Spark environment. It provides two dataset: a train data set and a test dataset. The train dataset is made of 878049 records and the test dataset, of 884262 records. Both of them contains incidents from January 1, 2003 to May 13, 2015. After pre-processing the data, the next step is to create and train the model.

The model would forecast the crime class using the weekday when the incident took place, the district where it happened and the time it occurred, as predictors. The following are features found in the crime dataset which is used for prediction of crime .

These are the attributes of the dataset:

| ATTRIBUTE | DESCRIPTION |
|---|---|
| Dates | Indicates the Timestamp of the incident. |
| Category | Category of the incident. |
| Descript | A short description of the incident. |
| DayOfWeek | Incident occurred on this Day of the week |
| PdDistrict | Under which police department the crime occurred |
| Resolution | Outcome of the incident. |
| Address | Address of the incident. |
| X | Longitude |
| Y | Latitude |

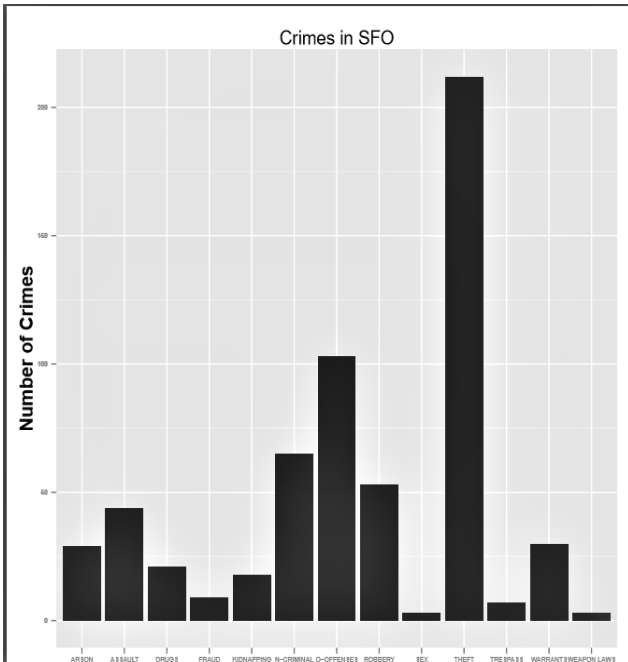The following fig 5.1 indicates the number of crimes happened during the respective years



**Fig :5.1 Classification of Crimes**
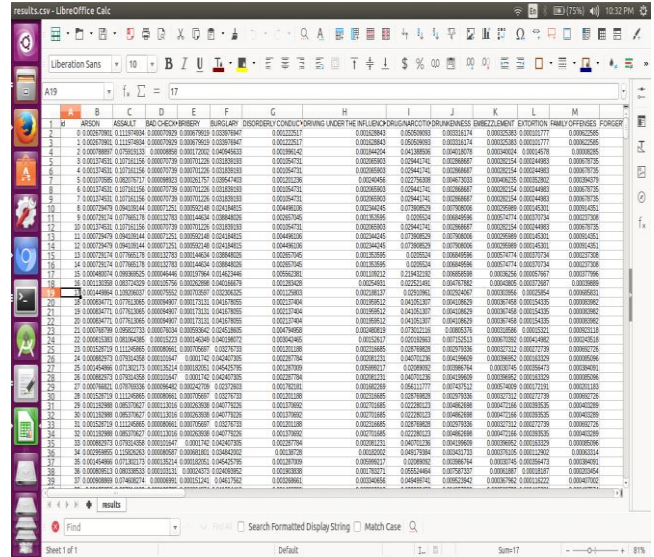


**Fig: 5.2 Multinomial Regression Values**



**Fig : 5.3 Prediction using MLR**

The fig 5.2 indicates the regression values obtained using multinominal regression techniques and fig 5.3 show the prediction of crime using MLR

## VI. CONCLUSION

The analysis of large crime data set is done by using R Studio and the Multinomial logistic regression algorithm. Apache Spark is used to achieve a high data processing. The crime factors change over time. Because of limited factors full accuracy cannot be achieved. More dependent crime attributes has to be added for getting better results in prediction. Advanced implementations like prediction by image processing can be implemented. Big Data Analytics is a real leap forward and a great opportunity to make huge gains in performance, productivity, income and profitability. The Big Data Era is here, and these are truly revolutionary times for business and technology experts to work together and delivering on the pledge.

## REFRENCES

1. Yerpude, P., & Gudur, V. (2017). Predictive modelling of crime dataset using data mining. *International Journal of Data Mining & Knowledge Management Process (IJDKP)*, 7(4).
2. Sathyadevan, S., & Gangadharan, S. (2014, August). Crime analysis and prediction using data mining. In *2014 First International Conference on Networks & Soft Computing (ICNSC2014)* (pp. 406-412). IEEE.
3. Bogomolov, A., Lepri, B., Staiano, J., Oliver, N., Pianesi, F., & Pentland, A. (2014, November). Once upon a crime: towards crime prediction from demographics and mobile data. In *Proceedings of the 16th international conference on multimodal interaction* (pp. 427-434). ACM.
4. A. Nasridinov, S. Ihm, and Y. Park. 'A Decision Tree Based Classification Model for Crime Prediction'. Information Technology Convergence, 531-538, 2013
5. David Smith., R Tops Data Mining Software Poll, Java Developers Journal, May 31, 2012.
6. Fox, John and Andersen, Robert (January2005). 'Using the R Statistical Computing Environment to Teach Social Statistics Courses'. Department of Sociology, McMaster University. Retrieved 2006-08-03.
7. SF Crime Classification Data- Kaggle. https://www.kaggle.com/sfcrime/data.
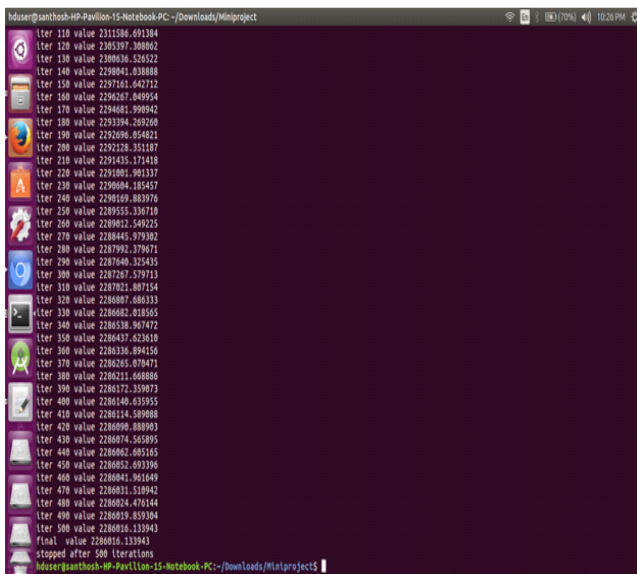8. YangchangZhoa(2012).R Data mining Examples and case studies. http://rdataminig.com.

9.   Feng, M., Zheng, J., Han, Y., Ren, J., & Liu, Q. (2018, July). Big Data Analytics and Mining for Crime Data Analysis, Visualization and Prediction. In International Conference on Brain Inspired Cognitive Systems (pp. 605-614). Springer, Cham.
10.  Das, P., & Das, A. K. (2019). Application of Classification Techniques for Prediction and Analysis of Crime in India. In Computational Intelligence in Data Mining (pp. 191-201). Springer, Singapore.
11.  Antolos, D., Liu, D., Ludu, A., & Vincenzi, D. (2013, July). Burglary crime analysis using logistic regression. In International Conference on Human Interface and the Management of Information (pp. 549-558). Springer, Berlin, Heidelberg.
12.  Kang, H. W., & Kang, H. B. (2017). Prediction of crime occurrence from multi-modal data using deep learning. PloS one, 12(4), e0176244

## AUTHORS PROFILE

**R. Raja Devi,** is received the B.E(Computer Science and Engineering) degree from Bharathiar University and also M.E(Computer Science and Engineering) degree from Annamalai University . She has 15 years of experience in teaching field. Her research interest includes web Technology , Data mining,optimization techniques,Machine learning , Service Oriented Architecture and Web Services Composition.Currently she is working as an Assistant Professor(Kongu Engineering College) and pursuing the Ph.D degree in Computer Science and Engineering. She has published four papers in international journals.

**E. M. Roopa Devi** is received the B.E(Electrical and Electronics) degree from Anna University and also M.E(Computer Science and Engineering) degree from Anna University . She has 9 years of experience in teaching field. Her research interest are in the area of Computer Networks, Network Security and Data mining.Currently she is working as an Assistant Professor(Kongu Engineering College) and pursuing the Ph.D degree in Computer Science and Engineering. She has published four papers in international journals.

**S.Vinoth Kumaris**, received the B.E(Electronics and Communication) degree from Anna University and also M.E(Computer and Communication) degree from Anna University .He has 6 years of experience in teaching field. His research interest are in the area of Computer Networks and Network Security.Currently ,He is working as an Assistant Professor in Kongu Engineering College.