

Minor Project Report Sem V - Vatsal, Ananya, Arjav (JIIT)

by Vatsal Gupta

General metrics

60,646

characters

10,957

words

726

sentences

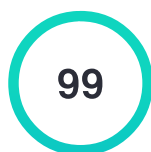
43 min 49 sec

reading
time

1 hr 24 min

speaking
time

Score



This text scores better than 99%
of all texts checked by Grammarly

Writing Issues

19

Issues left

12

Critical

7

Advanced

Plagiarism



31

sources

8% of your text matches 31 sources on the web
or in archives of academic publications

Writing Issues

13

Correctness

6

Misspelled words



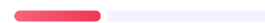
1

Improper formatting



2

Determiner use (a/an/the/this, etc.)



4

Confused words



6

Clarity

2

Wordy sentences



4

Passive voice misuse



Unique Words

15%

Measures vocabulary diversity by calculating the percentage of words used only once in your document

unique words

Rare Words

39%

Measures depth of vocabulary by identifying words that are not among the 5,000 most common English words.

rare words

Word Length

3.6

Measures average word length

characters per word

Sentence Length

15.1

Measures average sentence length

words per sentence

Minor Project Report Sem V - Vatsal, Ananya, Arjav (JIIT)

DATA MINING FOR SPATIO-TEMPORAL RULES WITH GEOSPATIAL ANALYSIS OF
CRIME PATTERNS
A PROJECT REPORT

Submitted by

VATSAL GUPTA [17104060]

ANANYA SHARMA [17104038]

ARJAV JAIN [17104070]

Under the guidance of

Mr. M. Gurve

(Assistant Professor, Department of Computer Science & IT) in partial
fulfillment for the award of the degree of

BACHELOR OF TECHNOLOGY

in

INFORMATION TECHNOLOGY

Sector 62, Noida, Uttar Pradesh 201309

NOVEMBER 2019

JAYPEE INSTITUTE OF INFORMATION TECHNOLOGY

20

(Deemed to be University under section 3 of UGC Act, 1956)

BONAFIDE CERTIFICATE

Certified that this project report titled "DATA MINING FOR SPATIO- TEMPORAL RULES WITH GEOSPATIAL ANALYSIS OF CRIME PAT-

TERNS " is the bonafide work of "VATSAL GUPTA [17104060], ANANYA SHARMA [17104038], ARJAV JAIN [17104070]", who carried out the project work under my supervision. Certified further, that to the best of my knowledge the work reported herein does not form any other project report or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

SIGNATURE

Mr. M. Gurve

SUPERVISOR

Assistant Professor

Dept. of Computer Science & IT

SIGNATURE

Ms. A. Kaur

PANEL MEMBER

Dept. of Computer Science & IT

□

SIGNATURE

Dr. V. Saxena

HEAD OF THE DEPARTMENT

Dept. of Computer Science & IT

SIGNATURE

Dr. A. Sahoo

PANEL MEMBER

Dept. of Computer Science & IT

ABSTRACT

Crime has always been a universal phenomenon. However, with the rapid urbanisation and industrialisation, there has been a surge in techno-industrial-urban complexes which offer a setting conducive to crime. Consequently, instances of crime have increased rapidly all across the globe among all segments of society. Furthermore, new forms of crime are emerging, and old forms are assuming new dimensions. Even though law enforcement agencies are doing their utmost to prevent crimes and catch criminals, they still face an uphill task. Comprehensive data is required to identify the criminals who are well organised and well equipped in using contemporary techniques that pose a severe threat to people's safety. Numerous investigations addressing this issue have generally employed disciplines of behaviour science and statistics. Recently, the data mining approach has proved to be a proactive decision-support tool in analysing, predicting and preventing crime. In this work, a framework based on clustering and association rule mining has been proposed to detect and analyse crime trend patterns from temporal and spatial crime

activity data. In addition, an open-source Geographic Information System (GIS) application, QGIS, is employed to reveal the overall crime hotspots as well as the hotspots of certain violent crimes in isolation. Furthermore, time series analysis of the crime has been done in order to analyse the changing patterns of crime with time. The resultant model can support police managers in developing more appropriate law enforcement strategies, as well as enhancing the use of police duty deployment for crime prevention.

PREFACE

This project aims to employ crime analysis techniques to enable the police administrators to objectively determine the nature of criminal activities in their jurisdictions and allow them to develop tactical action plans to combat them effectively. At the same time, this project can also help the public at large in being more aware and consequently safer. Doing this project has helped us to enhance our knowledge of data mining and data analyses techniques. Besides, this project has equipped us to learn a free and open-source GIS tool named QGIS, and, has contributed to our increased fluency in the python environment and its libraries.

ACKNOWLEDGEMENTS

²¹ | We would like to express our deepest gratitude to our mentor, Mr. Mahendra Gurve, for his valuable guidance, consistent encouragement, timely help, and providing us with an excellent atmosphere for doing our project. During the entire duration of our dissertation work, despite his busy schedule, he has extended cheerful and cordial support to us for completing this project work. We would also like to convey our sincere regards to all other faculty members of the department of CSE & IT, JIIT, who have bestowed their great effort and guidance at appropriate times. Our thanks and appreciations also go to people who have willingly helped in developing this project.

v

Contents

ABSTRACT

iv

PREFACE

v

ACKNOWLEDGEMENTS

vi

LIST OF TABLES

x

LIST OF FIGURES

xii

ABBREVIATIONS

1

1 INTRODUCTION

2

1.1 Problem Statement

2

1.2 Research Motivation

2

1.3 Research objectives

3

2 LITERATURE SURVEY

4

3 DATASET DESCRIPTION

8

3.1 Summary

9

4 RESEARCH METHODOLOGY

10

4.1 Clustering

10

4.1.1 K Means Clustering

10

4.2 Association Rule Mining

11

4.2.1 Apriori

12

4.2.2 Frequent Pattern (FP) Growth

12

4.3 QGIS

13

5 RESEARCH DESIGN

14

vi

5.1 Data Acquisition 14

5.2 Exploratory Data Analysis (EDA) 14

5.3 Data Pre-processing 17

5.4 Data Mining 18

5.4.1 K Means Clustering 18

5.4.2 Apriori 19

5.4.3 FP Growth	19
PROJECT SPECIFICATIONS 20	
6.1 Hardware Specifications	20
6.1.1 Recommended Configurations	20
6.1.2 Minimum Configurations	20
6.2 Software Specifications & Requirements	20
6.2.1 Project Compatibility	21
6.2.2 Python Packages	21
6.2.3 Python Modules	22
6.2.4 QGIS Plugins	22
ANALYSIS AND RESULTS 26	
7.1 K Means Clustering	26
7.2 Apriori	27
7.3 FP Growth	27
7.4 Hotspot Analysis	28
7.5 Cluster & Outlier Analysis	30
7.5.1 Interpretation of local Moran's I value	30
7.5.2 Interpretation of p-value	31
7.5.3 Interpretation of z-score	31
7.6 Time Series Analysis	32
EVALUATION METRICS 37	
Clustering Evaluation Metrics	37
8.1.1 Elbow Method	37
8.1.2 Silhouette Analysis	38
vii	
Association Rule Mining Evaluation Metrics	40
Comparison of FP Growth and Apriori	40

FUTURE DIRECTION AND CONCLUSION 43

9.1 Future Work 43

9.2 Conclusion 43

List of Tables

3.1 Dataset Columns 8

8.1 Comparison of Apriori and FP Growth 41

ix

List of Figures

3.1 CSV Dataset Snapshot	9
4.1 K Means Clustering	11
4.2 QGIS Desktop Snapshot	13
Number of crimes by day of the week	15
Number of crimes by month of the year	15
5.3 Number of crimes by type	16
Frequency of different types of crimes in different locations	16
Research Framework of our Project	17
Steps involved in K Means Clustering	18
5.7 Steps involved FP Growth	19
6.1 Step 1	23
6.2 Step 2	23
6.3 Step 3	24
6.4 Step 4	24
6.5 Step 5	24
6.6 Step 6 & 7	25
6.7 Step 8	25
6.8 Step 9	25
7.1 Clustering Results	26
7.2 Cluster Allotment	26
Apriori results with minimum support = 0.003	27
Apriori results with minimum support = 0.002	28
FP Growth results with minimum support = 0.003	29
Hotspot Analysis on the basis of count	29
Hotspot Analysis on the basis of density	30

x

Calculation of local Moran's I value, p-value and z-score	31
Cluster & Outlier Analysis on Chicago Dataset	32
7.10 Heatmap - January 2017	33
7.11 Heatmap - February 2017	33
7.12 Heatmap - March 2017	33
7.13 Heatmap - April 2017	34
7.14 Heatmap - May 2017	34
7.15 Heatmap - June 2017	34
7.16 Heatmap - July 2017	35
7.17 Heatmap - August 2017	35
7.18 Heatmap - September 2017	35
7.19 Heatmap - October 2017	36
7.20 Heatmap - November 2017	36
7.21 Heatmap - December 2017	36
8.1 Computing WCSS	38
8.2 The Elbow method	38
8.3 Silhouette Analysis	39
8.4 Time Taken using Apriori	42
8.5 Time Taken using FP Growth	42

xi

ABBREVIATIONS

GIS Geographic Information System

FP Frequent Pattern

HDD Hard Disk Drive

SSD Solid State Drive

IUCR Illinois Uniform Crime Reporting EDA Exploratory Data Analysis WCSS

Within Cluster Sum of Squares SQL Structured Query Language

TAR Temporal Association Rules

SVM Support Vector Machines

KNN K Nearest Neighbour

xii

Chapter 1

INTRODUCTION

Problem Statement

22

Increased population, technological advancements and heightened competition for economic resources have created various social problems. Many of these changes in the human condition have brought new challenges to the doorstep of the law enforcement profession that begs for resolution. The major challenge facing law enforcement agencies is to deal with the increased number of criminal activities effectively and efficiently. Current policing strategies work towards finding the criminals, basically after the crime has occurred. However, with the help of technological advancement, we can use historical crime data to recognise crime patterns [8]. If enforcement agencies have a prior assumption of the class of the crime, it would give them tactical advantages and help resolve cases faster. An overall study of criminal activity in a geographic area also helps in understanding the underlying pattern of the crime in that area.

Research Motivation

Criminals have been a nuisance for society in all corners of the world for a long time now. Measures are required to eradicate crimes from all over the world or at least limit its occurrence. The large volumes of crime datasets, as well as the complexity of associations between these kinds of data, have made criminology an appropriate field for the application of data mining techniques. Criminology is an area that focuses on the scientific study of crime, criminal behaviour, and law enforcement. In simpler terms, it is a process that aims to identify crime characteristics [10]. It is one of the most relevant fields where

the application of data mining techniques can produce remarkable results that can help and support police forces. Using data mining and data analytic techniques, we can analyze the crime patterns that can further help the authorities to understand the underlying reasons for the occurrence of crime and can, therefore, help them to a large extent in the prevention of future crimes.

1

Research objectives

The primary objective of our work is to analyse criminal data based on demographics, spatial and temporal information and consequently identify useful crime patterns to aid police in preventing crimes. Towards this end, we have employed data mining and crime mapping techniques. The main objectives of our project work are summarised as follows:

Identifying the crime patterns based on a criminal dataset that contains the geographical location and basic details of the criminal activity.

Exploring data mining techniques to generate association rules for crime analysis.

Visualising these patterns on an open source GIS software - QGIS for better understanding of the results.

Chapter 2

LITERATURE SURVEY

Much work has been done in the direction of crime analysis to improve the activities aimed at detecting and preventing safety problems for the public. Techniques ranging from conventional data association methods to the modern approaches of data mining have been applied to this field. This section aims to summarise the work done in this regard.

In [2], automated approaches to data association to increase the accuracy of crime prediction have been proposed. Results included in the paper indicated that the employed data association methods significantly reduced the time required by manual methods while maintaining a high level of accuracy, comparable to that of experienced crime analysts. Furthermore, in contrast to existing analysis techniques that employed Structured Query Language (SQL), these methods were both faster and more accurate.

Nath (2006) proposed a clustering technique over supervised learning techniques such as classification for crime data analysis [15]. In this work, K Means clustering has been applied to identify criminal patterns and subsequently to help prevent future crimes. Furthermore, the clustering algorithm has been integrated with a geospatial plot using which a crime analyst can choose a time range and one or more types of crime from specific geography and view the result graphically.

The concept of Temporal Association Rules (TAR) was introduced in [16] to solve the problem of time series handling by including time expressions into association rules. An incremental algorithm - ITAR (Incremental TAR) has been

proposed in this paper to overcome the re-scanning issue in the existing TAR algorithm for updating the dataset. This paper made use of negative border method for preserving temporal association rules with numerical attributes. The temporal negative border method proposed in this paper only retains all the past winners who become losers in subsequent rounds instead of maintaining a power set of dense base cubes. As a result, the number of losers of base cubes held by the negative temporal boundary was minimised. Preliminary results showed a significant improvement over the recurrent TAR algorithm and showed that ITAR is very stable and has good performance.

3

[11] examined the performance of two types of Support Vector Machines (SVM) techniques: two-class SVMs and one-class SVMs for predicting the location of crime hotspots when a predefined level of crime rate is given. The paper also compared an SVM with a neural network-based approach and a spatial auto-regression-based approach. Initially, K means clustering has been used as a data selection approach. After labelling the data, the resulting dataset has been used as the input of the SVM algorithm. In both one-class and two-class SVMs, different kernel functions were used to determine the accuracy of the classification. Their experiments have shown that one-class SVM produces reasonable results, particularly when the size of the training set is small with more positive samples. However, for larger datasets and more negative samples, two-class SVMs were found to have a better performance.

23

A framework of intelligent decision-support model based on a fuzzy self-organising map (FSOM) network to detect and analyse crime trend patterns from temporal crime activity data was proposed in [13]. Besides, a rule

extraction algorithm has been employed to uncover hidden causal-effect knowledge and reveal the shift around effect. As per the analysis of the experimental results, they discovered characteristics of four crime patterns, namely, typical, gradual increase, sharp increase, and Wintertime. Their results showed that their framework could help provide vital information to the police management for determining the kind of duty deployment that should be employed. However, a significant limitation of their study was the lack of evaluation metrics to evaluate the accuracy of their model.

[18] discussed the preliminary results of a crime forecasting model developed in collaboration with the police department of a United States city in the Northeast. Their datasets comprised of aggregated counts of crime and crime-related events categorised by the police department. The location and time of these events were also included in the dataset. This work employed various data classification techniques to perform crime forecasting and to determine the best classification method for predicting crime hotspots. Their results indicated that 1NN classifier modified with location constraint is better in finding similar circumstances in a neighbourhood but not in the entire city. Furthermore, their results found Naïve Bayes classifier to be the best probability predictor.

In [7], two classification methods, namely, Naïve Bayes and Decision Tree, have been applied to a crime dataset to predict 'Crime Category' for different states of the United States of America. The results obtained from the experiment indicated that the Decision Tree algorithm

outperformed the Naïve Bayesian algorithm and achieved 83.9519% Accuracy in predicting 'Crime Category' for different states of the US.

Various strategies to find spatial and temporal criminal hotspots have been employed in [1]. Freely available crime datasets of two cities, namely, Denver (Colorado) and Los Angeles (California) have been used in this paper. The crime patterns in this paper were achieved by applying the Apriori algorithm on both the Denver and the Los Angeles datasets based on a predefined threshold. 62 interesting frequent patterns were generated for Denver, whereas 59 interesting frequent patterns were generated for Los Angeles. Two classification methods, namely, Naïve Bayes and Decision Tree classifiers, were also employed for crime type prediction. Their results indicated that Naïve Bayes classifier had better accuracy (51%) as compared to the Decision Tree classifier (42%). Furthermore, this paper presented a statistical analysis of the dataset supported by several graphs.

Spatio-temporal and demographic data has been used in [3] to predict which category of crime is most likely to have occurred at a given time and place. Various classification algorithms such as Naive Bayes, Support Vector Machines, Gradient Boosted Decision Trees, and Random Forests have been applied and compared. The inputs to these algorithms were time (hour, day, month, year), place (latitude, longitude, and police district), and demographic data (population, median income, minority population, and the number of families). The output produces was the category of the crime that is likely to have occurred. Classification of blue- and white-collar crimes, as well as violent and non-violent crimes, was also done using the algorithms as mentioned above. The Gradient Boosted Decision Trees algorithm was found to have the best accuracy.

Linear Regression, Additive Regression, and Decision Stump algorithms were implemented in [14] on the same finite set of features. The results indicated that the linear regression algorithm has the best performance among the three selected algorithms. The relatively poor performance of the Decision Stump algorithm was attributed to a certain factor of randomness in the various crimes and the associated features; the branches of the decision trees are more stringent and give accurate results only if the test set follows the pattern modelled. On the other hand, the linear regression algorithm was able to handle the randomness in the test samples to a great extent. This paper explored the efficiency and accuracy of the machine learning algorithms in data mining research for predicting trends of violent crimes.

5

In [17], a model for estimating the regions with high probability of crime incidence has been proposed. Various data mining techniques, in particular, data clustering algorithms have been employed to extract unknown but useful information from unstructured data. Various clustering techniques like K - Means and Fuzzy C have been used in this paper to take into account the dynamic nature of the crimes.

In [12], Vancouver crime dataset for the last 15 years was used for the application of predictive machine learning models such as K Nearest Neighbour (KNN) and enhanced decision trees were used and gave crime prediction accuracies between 39% and 44%. Owing to the use of different approaches, performance, complexity, and training time of algorithms were

found to be different. Although the predictive model used in this paper lacks reliability, it can act as a basis for further studies.

A novel methodology for forecasting crime has been applied in [9]. KNN has been used in this paper to calculate optimal values for good performance. Furthermore, a Bayesian Net- work was used to establish associations that benefit a range of variables. The results of the simulations indicated that the Naïve Bayes algorithm was highly precise and time-efficient.

6

Chapter 3

DATASET DESCRIPTION

The dataset used for the purposes of this work has been retrieved from the official government portal of the City of Chicago. The dataset reflects reported incidents of crime (except for murders) that occurred in the City of Chicago from 2001 to November of 2019. Addresses are manifested at the block level only, and specific locations are not identified to preserve the privacy of crime victims. Although the original dataset has 6.6 million rows, owing to the computational constraints, we have chosen a small fraction of this dataset. The description of all the attributes(columns) is given in the table below:

Column Name

Description

Data Type

ID

Unique Identifier for the record

Number

Case Number

27

The Chicago Police Department RD Number (Records Division Number), which is unique to the incident.

Plain Text

Date

Date when the incident occurred.

Date & Time

Block

The partially redacted address where the incident occurred, placing it on the same block as the actual address.

Plain Text

IUCR

The Illinois Uniform Crime Reporting code. This is directly linked to the Primary Type and Description.

Plain Text

Primary Type

The primary description of the IUCR code.

Plain Text

Description

The secondary description of the IUCR code, a subcategory of the primary description.

Plain Text

Location Description

Description of the location where the incident occurred.

Plain Text

Arrest

Indicates whether an arrest was made.

Checkbox

Domestic

Indicates whether the incident was domestic-related as defined by the Illinois

Domestic Violence Act

Checkbox

Beat

Indicates the beat where the incident occurred.

Plain Text

District

Indicates the police district where the incident occurred.

Plain Text

Ward

The ward (City Council district) where the incident occurred.

Number

Community Area

Indicates the community area where the incident occurred.

Plain Text

FBI Code

Indicates the crime classification as outlined in the FBI's National Incident-Based Reporting System (NIBRS).

28

Plain Text

X Coordinate

The x coordinate of the location where the incident occurred.

Number

Y Coordinate

The y coordinate of the location where the incident occurred

Number

Year

Year the Incident occurred

Number

Updated On

Date and Time the record was last updated

Date & Time

Latitude

The latitude of the location where the incident occurred.

Number

Longitude

The longitude of the location where the incident occurred.

Number

Location

The location where the incident occurred in a format that allows for creation of maps and other geographic operations on this data portal.

Location

Table 3.1: Dataset Columns

Additional information about some attributes -

IUCR: Illinois Uniform Crime Reporting (IUCR) codes are four-digit codes that are used by law enforcement agencies to classify criminal incidents when taking individual reports. In addition, these codes are used to aggregate types of cases for statistical purposes.

The Chicago Police Department currently uses more than 350 IUCR codes to classify criminal offences. The list of IUCR codes is available at <https://data.cityofchicago.org/d/c7ck-438e>.

7

Beat: The smallest police geographic area is known as a beat – each beat has a dedicated police beat car. Three to five beats constitute a police sector, and three sectors further constitute a police district. There are a total of 22 police districts in the Chicago Police Department.

Community Areas: Chicago has a total of 77 community areas.

Coordinates: All the coordinates, including latitude and longitude, are projected in State Plane Illinois East NAD 1983 projection. Furthermore, the coordinates (including latitude and longitude) map to a location that is shifted from the actual location for partial redaction but falls on the same block to protect the privacy of the crime victim.

Summary

Dataset Name: Crimes - 2001 to present

Number of Rows: 6.6 million

Number of Columns = 22

Dataset Source - Chicago Police Department's CLEAR (Citizen Law Enforcement Analysis and Reporting) system -

<https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-present/ijzp-q8t2>

Dataset Availability - Publicly Available for free

Figure 3.1: CSV Dataset Snapshot

8

Chapter 4

RESEARCH METHODOLOGY

Data mining techniques include supervised learning methods such as regression and classification algorithms and unsupervised learning methods such as clustering algorithms and association rule mining algorithms. In this work, to analyse the data and observe the crime patterns, three unsupervised learning algorithms, namely, 1) K Means (Clustering), 2) Apriori and, 3) FP Growth (Association Rule Mining) have been used. These algorithms have been applied to a segment of the crime dataset available freely on the official website of the Chicago Police Department. Furthermore, to validate the results obtained from the application of the algorithms as mentioned above, an open-source GIS software - QGIS, has been used to map the crime dataset. In addition to this, various graphical analysis techniques such as hotspot analysis, cluster & outlier analysis and, time series analysis have been performed in order to extract meaningful patterns and characteristics from the crime data.

All of the algorithms and techniques, as mentioned above, have been further elaborated in the following sections.

Clustering

Clustering is an unsupervised data mining strategy to group the relevant data into desired clusters. The centroid of each cluster represents the collection of feature values that define the resulting groups. By analysing in which cluster the data points lie, clustering enables us to gain valuable insights about our data. A clustering technique has been chosen over any other supervised technique such as classification since crimes widely vary in nature, and crime databases are often filled with unsolved crimes. Therefore, a classification technique that will rely only on the existing and known solved crimes will not give good predictive quality for future crimes.

K Means Clustering

K-means clustering is the simplest and most commonly used clustering algorithm owing to its less computational complexity. It is used when the data available is unlabelled (i.e., data

9

30 without defined categories or groups) or when the data sets are large. The algorithm works iteratively to assign each data point to one of the K groups

based on the features that are provided. Data points are clustered based on feature similarity.

Figure 4.1: K Means Clustering

Association Rule Mining

Association rule mining is a technique for recognizing frequent patterns and associations among a set of objects. Association rules are used to help uncover relationships between seemingly unrelated data in a set of transactions. An association rule has two components, an antecedent or the LHS (if) and a consequent or the RHS(then). An antecedent is an item that exists in the list of transactions. A consequent is an item that is found in combination with the antecedent. Consider A and B to be two items present in a list of N transactions, where A is the antecedent and B is the consequent, then:

Support: Support of item A is the ratio of the number of transactions in which item A appears to the total number of transactions.

Support (A) =

$\frac{\text{frequency (A)}}{N}$

(4.1)

Confidence: Confidence measures the percentage of times item B is purchased, given that item A was purchased. Values of confidence range from 0 to 1, where

0 indicates that B is never purchased when A is purchased, and 1 indicates that B is always purchased whenever A is purchased.

Confidence ($A \rightarrow B$) =

10

$\frac{\text{frequency}(A, B)}{\text{frequency}(A)}$ (4.2)

frequency (A)

Lift: The lift of a rule measures the percentage of times item B is purchased when item A is purchased while controlling the popularity of item A. A lift value of 1 implies that there is no association between the items.

Support (A, B)

$\text{Lift}(A \rightarrow B) = \frac{\text{Support}(A, B)}{\text{Support}(A) * \text{Support}(B)}$ (4.3)

Conviction: It represents by what factor the correctness of the rule would reduce if the antecedent (in this case - A) and the consequent (in this case - B) of the rule were independent. Higher is the confidence; higher is the conviction of a rule.

\rightarrow

$\text{Conviction}(A \rightarrow B) = \frac{1 - \text{Support}(B)}{1 - \text{Confidence}(A \rightarrow B)}$

1 – Confidence ($A \rightarrow B$)

\square

(4.4)

Apriori

31

Apriori is the most commonly used association rule mining algorithm to find frequent itemsets. It takes advantage of the fact that any subset of a frequent itemset is also a frequent itemset. The algorithm, thereby, reduces the number of candidates being considered by only exploring the item-sets whose support count is greater than the minimum support count. Steps to implement the Apriori algorithm are as follows:

Set minimum support and confidence.

Take all the subsets in transactions having higher support than minimum support.

Take all the rules of the subsets having higher confidence than minimum confidence.

Sort rules by decreasing lift.

FP Growth

FP Growth is a tree based association rule mining algorithm. Unlike Apriori, FP Growth uses a pattern growth approach rather than a candidate generation approach to find frequent itemsets. FP Growth is faster and more efficient as compared to Apriori as it scans the dataset only twice, whereas Apriori scans the dataset multiple times in order to generate frequent itemsets.

11

QGIS

A GIS is a system designed for the collection, storage, manipulation, analysis, management, presentation and analysis of spatial or geographical data. GIS applications allow users to analyze geospatial information, edit data in maps, and present the results of all these operations. QGIS is one such free and open source cross-platform desktop application. QGIS supports shapefiles, coverages, personal geodatabases, dxf, MapInfo, PostGIS, and other formats. It also supports Web services, including Web Map Service and Web Feature Service, to allow the use of data from external sources. QGIS has a plethora of plugins that enable a dynamic analysis of data including hotspot analysis, time series analysis, cluster & outlier analysis.

Figure 4.2: QGIS Desktop Snapshot

12

Chapter 5

RESEARCH DESIGN

Data Acquisition

The dataset under scrutiny has been acquired from the official website of the Chicago City Government. The dataset contains the information of all the crime incidents from 2001 to 2019 with the exception of murders. A detailed description of the dataset has already been presented in Chapter 3.

EDA

EDA is a systematic way of visualisation and transformation to explore and summarise the main characteristics of the available data. The main objectives of an EDA include:

Generating questions about the available data.

Searching for answers by visualising, transforming, and modelling the data.

Using the findings to generate new questions.

In our project, the EDA was carried out on the crime data for the current year, i.e., 2019.

The following criteria were considered while performing the EDA:

Number of crimes by day of the week. (Figure 5.1)

Number of crimes by month of the year. (Figure 5.2)

Number of crimes by type. (Figure 5.3)

Frequency of different types of crime in different locations. (Figure 5.4)

13

Figure 5.1: Number of crimes by day of the week

Figure 5.2: Number of crimes by month of the year

Note: November's crime count is less due to the availability of data only until the first week of November.

14

Figure 5.3: Number of crimes by type

Figure 5.4: Frequency of different types of crimes in different locations

15

Data Pre-processing

Data pre-processing is a data mining procedure that involves converting raw data into an understandable format. Real-world data is often incomplete, inconsistent and lacking in certain behaviours or trends and is likely to contain several errors. Data pre-processing is a proven method of addressing these issues. Data Pre-processing steps involved in our project:

Dropping redundant columns such as 'X Coordinate', 'Y Coordinate', Location (since 'Latitude' and 'Longitude' are available), 'Updated On' and more.

Converting the date and time given in the dataset to pandas DateTime format.

For clustering, it was necessary to normalise the time, district and the IUCR code values in order to form clusters based on crime type, time and location.

The formula adopted for normalisation is given below:

$$x = x - x_{\text{minimum}}$$
$$\square$$

(5.1)

normalised

 $\square x_{\text{maximum}}$ $\square x_{\text{minimum}}$

Figure 5.5: Research Framework of our Project

16

Data Mining

K Means Clustering

The steps involved while performing K Means clustering on the Chicago crime dataset are given in the figure 5.6.

Figure 5.6: Steps involved in K Means Clustering

17

Apriori

The steps involved in applying the Apriori algorithm have been discussed in Chapter 4.

FP Growth

The flowchart below (figure 5.7) reveals the steps involved in generating association rules from our dataset.

Figure 5.7: Steps involved FP Growth

18

Chapter 6

PROJECT SPECIFICATIONS

Hardware Specifications

Recommended Configurations

CPU Speed: 1.4 GHz or higher recommended (per core)

Processor: Intel Core i5-6100 (6th Gen) Dual Core or above

Memory/ RAM: 8GB or higher

Monitor/ Display: 13" LCD monitor, resolution of 1920 x 1080 or better

Storage: 200 GB Hard Disk Drive (HDD)/ 128 GB Solid State Drive (SSD) or above

Minimum Configurations

CPU Speed: 1.2 GHz (per core)

Processor: Intel Core i3-6100 (6th Gen) Dual Core

Memory/ RAM: 4GB

Monitor/ Display: 13" LCD monitor, resolution of 1600 x 900

Storage: 100 GB HDD/ 64 GB SSD

Software Specifications & Requirements

Operating System Used: Microsoft Windows 10 Version 10.0.17763 64 bit

Programming Language Used: Python 3.7.3 64 bit

-

IDEs/ Environments Used:

Spyder 3.3.6

VS Code 1.40.1

Jupyter Notebook 6.0.2

GIS Application Used: QGIS Desktop 3.10.0 A Coruña based on Python 3.7.0

19

Project Compatibility

-

Operating Systems:

Microsoft Windows 7 or above Mac OS El Capitan or above

Linux Flavours - Debian/Ubuntu, Fedora, Mandriva, Slackware, ArchLinux,

Flatpak, openSUSE, RHEL

Python Version: Python 3.x.x

-

IDEs: Compatible with all Python IDEs provided all the required packages and modules are installed

Recommended IDEs: Spdyer 3.3.x and VS Code 1.33.1 or above

Python Packages

pandas

numpy

matplotlib

efficient_apriori

tkinter

lpython

sklearn

mpl_toolkits

plotly

pyfpgrowth

yellowbrick

pysal (for QGIS)

rtree (for QGIS)

gdal (for QGIS)

geopandas (for QGIS)

20

Python Modules

time

datetime

apriori (from efficient_apriori)

future

pylab

Axes3D (from mpl_toolkits)

webbrowser

warnings

SilhouetteVisualizer (from yellowbrick)

fp (from pyfpgrowth)

esda (from pysal)

spreg (from pysal)

QGIS Plugins

TimeManager

HotspotAnalysis

Qgis2threejs

MapSwipe Tool

How to install the modules and packages in Python?

In case of standalone python IDE:

Using pip3 installer in terminal.

In case of anaconda environment:

Using conda install or pip3 installer in anaconda prompt.

21

How to install the plugins in QGIS?

Download plugins in zip format from <https://plugins.qgis.org/>

Open QGIS Desktop Application and open the 'Plugins' tab

Choose 'Manage and Install Plugins'

A new window would open, click on 'Abort Fetching'

After the completion of the previous step, another window would open. Go To
'Install from ZIP'

Select your zip file

Click on 'Open'

Click on 'Install Plugin'

You can verify that your plugin has been installed in the 'Installed' option.

Figure 6.1: Step 1

Figure 6.2: Step 2

22

Figure 6.3: Step 3

Figure 6.4: Step 4

Figure 6.5: Step 5

23

Figure 6.6: Step 6 & 7

Figure 6.7: Step 8

Figure 6.8: Step 9

24

Chapter 7

ANALYSIS AND RESULTS

K Means Clustering

Figure 7.1 shows the clustering results obtained using K Means Clustering.

Clustering with K=4 (b) Clustering with K=5

Figure 7.1: Clustering Results

Figure 7.2 depicts the cluster allotment. Here, "Normalized_time" denotes the value of time between 0 and 1. Lower values in this column would indicate midnight to early morning, medium

Figure 7.2: Cluster Allotment

25

values would indicate the afternoon sessions, and high values would indicate the evening and night time.

Apriori

For a minimum support of 0.003 and minimum confidence of 0.005, Apriori generated only two rules (refer figure 7.3), whereas for a minimum support of 0.002 and the same minimum confidence, Apriori generated 70 rules (refer figure 7.4). The antecedent (LHS) of the rules generated using both, Apriori and FP Growth, have 3 or more of the following attributes: ["Block", "Location Description", "District", "Timeslot", "Month"], whereas the consequent (RHS) is always "Primary Type".

Figure 7.3: Apriori results with minimum support = 0.003

FP Growth

In contrast to Apriori, FP Growth generated a large number even for relatively higher values of minimum support. However, to compare the performance of these two algorithms, the minimum support and confidence were kept the same. For the minimum confidence value of 0.005,

26

Figure 7.4: Apriori results with minimum support = 0.002

a total of 450 rules were generated when the value of minimum support was kept 0.003 (refer figure 7.5).

Hotspot Analysis

An additional plugin named "Hotspot Analysis" has to be downloaded using the steps shown in figures 6.1 to 6.8 to perform hotspot analysis in QGIS.

Furthermore, additional python packages, as mentioned in Chapter 6 (Section 6.2.2 - Python Packages), need to be installed.

Hotspot analysis is a spatial analysis and mapping method interested in the identification of clustering of spatial events. These spatial events are depicted as points in a map and refer to locations of events or objects. In our work, hotspot analysis has been performed to identify locations with high occurrences of crime. Hotspot analysis has been carried out in our project to identify areas of high crime incidents. Hotspot analysis was carried out on the basis of both count (figure 7.6) and concentration (figure 7.7) of crime (crime count/area of the location). By carefully looking at the figures 7.6 and 7.7, it can be observed that certain census tracts were classified as areas of low crime counts in figure 7.6, even though they were classified as areas of high crime density. This is because, although these areas have less crime count compared to other census tracts, they have high crime count to area ratio due to their small areas.

Figure 7.5: FP Growth results with minimum support = 0.003

Figure 7.6: Hotspot Analysis on the basis of count

28

Figure 7.7: Hotspot Analysis on the basis of density

Cluster & Outlier Analysis

Cluster & Outlier Analysis is performed to identify spatial clusters of features with high or low values as well as spatial outliers. To do this, a Cluster & Outlier tool that calculates a local Moran's I value, a z-score, a pseudo p-value is used. [5] The cluster/outlier type (COType) field differentiates between a statistically significant cluster of low values (LL), cluster of high values (HH), outlier in which a low value is surrounded mainly by high values (LH), and an outlier in which a high value is surrounded mostly by low values (HL).

Interpretation of local Moran's I value

A positive I value indicates that a feature has neighbouring features with similar values (either high or low), i.e., this feature is part of a cluster. A negative I value indicates that a feature has neighbouring features with different values, i.e., this feature is an outlier.

29

Figure 7.8: Calculation of local Moran's I value, p-value and z-score

Interpretation of p-value

The p-value measures the possibility that the observed spatial pattern was created by some random process. A low p-value signifies the unlikeliness of the spatial pattern being generated by some random process.

Interpretation of z-score

-

32 | A high positive z-score for a feature signifies that the surrounding features have similar values (either high or low). The COType field in the Output Feature Class will be LL for a statistically significant cluster of low values and HH for a statistically significant cluster of high values respectively. [6]

A low negative z-score for a feature signifies a statistically significant spatial data outlier.

30

Depending on the values of the surrounding features, The COType field in the Output Feature Class will either be LH (low feature value surrounded by features with high values) or HL (high feature value surrounded by features with low values).

Figure 7.9: Cluster & Outlier Analysis on Chicago Dataset

Time Series Analysis

Time series analysis comprises techniques for analysing time series data in order to extract meaningful statistics and other characteristics of the data. In our project, we have performed a time series analysis using QGIS. Instead of forecasting crimes for the next time-period, we have limited our project to just the analysis of previous trends.

31

Figure 7.10: Heatmap - January 2017

Figure 7.11: Heatmap - February 2017

Figure 7.12: Heatmap - March 2017

32

Figure 7.13: Heatmap - April 2017

Figure 7.14: Heatmap - May 2017

Figure 7.15: Heatmap - June 2017

33

Figure 7.16: Heatmap - July 2017

Figure 7.17: Heatmap - August 2017

Figure 7.18: Heatmap - September 2017

34

Figure 7.19: Heatmap - October 2017

Figure 7.20: Heatmap - November 2017

Figure 7.21: Heatmap - December 2017

35

Chapter 8

EVALUATION METRICS

Evaluation metrics are used to verify a model's performance. It is crucial to choose proper metrics in order to evaluate how well an algorithm is performing. The evaluation metrics used for different models are briefly explained in the following sections:

Clustering Evaluation Metrics

In comparison to the supervised learning methods, where we have the ground truth to evaluate the model's performance, clustering analysis does not have a robust evaluation metric that we can use to evaluate the outcome of different clustering algorithms. Moreover, since K-Means requires K as input and does not learn it from data, there is no right answer in terms of the number of clusters that we should have in any problem. Sometimes domain knowledge and intuition may help, but usually, that is not the case [4]. In the cluster-predict methodology, we can evaluate how well the models are performing based on different K clusters since clusters are used in the downstream modelling. In our work, we have used two such metrics that may give us better intuition about K:

Elbow Method

Silhouette Analysis

Elbow Method

35 | The elbow method is one of the most common and technically robust methods for determining a good value of K for clustering. It provides us with an idea of what a good number of clusters would be on the basis of the Within Cluster Sum of Squares (WCSS) between data points and their assigned clusters' centroids. For $K = 3$, WCSS can be computed as follows:

36

Figure 8.1: Computing WCSS

While clustering performance as measured by WCSS might increase (i.e. WCSS decreases) with an increase in K , the rate of increase usually decreases. So, we choose that that value of K where WCSS starts to flatten out and to form an elbow, i.e., where the rate of increase in performance starts to diminish.

Figure 8.2: The Elbow method

The graph above shows that $K = 4$ and $K = 5$ are excellent choices for the value of K for our dataset. Although the elbow method provides useful insights for choosing the value of K , sometimes it is difficult to figure out a good number of clusters to use by looking at the elbow curve because the curve is monotonically decreasing and may not show a distinct point where the curve

starts flattening out. In such cases, silhouette analysis can be of significant help.

Silhouette Analysis

36

Silhouette analysis is employed to interpret and validate the consistency within clusters of data [?]. The silhouette value is a measure of how similar an entity is to its cluster (cohesion) compared to other clusters (separation). It can be used to determine the degree of separation between clusters. For each sample, the steps involved in silhouette analysis are as follows:

The mean distance from all data points in the same cluster (a_i) is calculated.

The mean distance from all data points in the closest cluster (b_i) is computed.

37

The silhouette coefficient is computed using the equation given below:

–

 $b_i a_i$
$$\text{Coefficient} = \max(a_i, b_i)$$

□

(8.1)

37

The coefficient can take values in the interval $[-1, 1]$:

If it is 0 \rightarrow the sample is very close to the neighbouring clusters.

If it is 1 \rightarrow the sample is far away from the neighbouring clusters.

If it is -1 \rightarrow the sample is assigned to the wrong clusters.

With 3 Clusters (b) With 4 Clusters

(c) With 5 Clusters (d) With 6 Clusters

Figure 8.3: Silhouette Analysis

For a good value for the number of clusters, the following conditions should be satisfied:

The mean value of silhouette coefficients should be as close to 1 as possible.

The silhouette plot of each cluster should be above the mean silhouette value as much as possible. Any plot region below the mean value is not desirable.

The width of the plot should be as uniform as possible.

38

From the fig. 8.3, it can be established that $K = 4$ and $K = 5$ are the optimal values for K . For $K = 3$ and $K = 6$, not only is the mean silhouette value low, the

silhouette plot of many individual clusters is below that the average silhouette line.

For our dataset, Silhouette analysis has further corroborated the results obtained from the el- bow method. This demonstrates the effectiveness of our implementation of K Means clustering algorithm.

Association Rule Mining Evaluation Metrics

Like clustering, association rule mining algorithms lack properly defined evaluation metrics. However, support, confidence, lift, and conviction of the rules, as established in eqns. 4.1, 4.2, 4.3 and 4.4 respectively, can offer insights pertaining to the correctness of the rules. The evaluation metrics for association rule mining algorithms can broadly be classified into two categories:

Symmetric Measures: Support, Lift/Interest etc.

Asymmetric Measures: Confidence, Conviction etc.

Comparison of FP Growth and Apriori

Since FP Growth and Apriori have been used for the same purpose, i.e., extracting rules concerning crime characteristics from our dataset, it becomes imperative to compare their performance. We have evaluated and compared the performance of these two algorithms on the basis of two parameters:

The values of lift and conviction (for the same values of minimum support and confidence).

Time and space complexity.

Lift and Conviction

Conviction of a rule measures by what factor the correctness of the rule would reduce if the antecedent and the consequent of the rule were independent. For example, if the conviction

39

value of a rule is 1.2, then, the rule would be incorrect 20% more often if the association between antecedent and consequent was purely by chance. Some interesting points to note:

If items are independent: $\text{lift} = \text{conviction} = 1$

-

Higher is the value of confidence, higher is the value of conviction. If confidence of a rule is equal to 1, then conviction of that rule would be infinity.

From the rules generated using Apriori and FP Growth, the average conviction of the rules obtained using FP Growth was found to be higher than the average conviction values of the rules obtained using Apriori.

Time & Space Complexity

FP Growth makes only two passes over the dataset as compared to Apriori, which scans the dataset multiple times. Reduced number of scans make FP Growth a more time-efficient algorithm. Furthermore, FP Growth uses a tree data structure for storing the generated rules, which subsequently make it more space-efficient as well.

Table 8.1: Comparison of Apriori and FP Growth

Apriori

FP Growth

Time Complexity

$O(2d)$

$O(n+d)$

Space Complexity

$O(2d)$

$O(nd)$

Table 8.1 compares the time and space complexity of the two algorithms. Here, n is the number of transactions, i.e., rows in our dataset, while d denotes the number of unique items in our dataset.

To determine the time taken by both the algorithms, namely, FP Growth and Apriori for rule generation, a timing module was used¹. For the same dataset, it was observed² that FP Growth was significantly faster than Apriori in generating rules. Apriori took around 6.28 seconds to generate rules (refer figure 8.4), whereas, FP Growth took only 0.25 seconds (refer figure 8.5). This discrepancy in the time taken would be more substantial for larger datasets.

40

Figure 8.4: Time Taken using Apriori

Figure 8.5: Time Taken using FP Growth

41

Chapter 9

FUTURE DIRECTION AND CONCLUSION

Future Work

Taking a close look at the problems discussed, we can say that our model is pretty much expandable to other domains as well. However, there exist certain limitations in the proposed model. One limitation of our model is its sensitive nature to the quality of input data that may be inaccurate or have missing information. Another limitation of the current work is the difficulty of evaluating the accurate performance of the proposed model, so future works could be aimed³ at finding a way to use quantified methods to evaluate the degree of crime reduction achieved, the improvement in duty deployment, and

the impact on society. Moreover, owing to computational constraints, algorithms have only been applied to data for the current year. If the period under consideration could be increased⁴, the results produced would be more reliable.

As a future extension of our work, supervised learning models such as classification models can be used to predict the type of crime. It is also a useful extension for our study to consider neighbourhood income information in order to see if there is a relationship between the income level of neighbourhoods and their crime rate. Additionally, the ability to search suspect description in regional, FBI databases, traffic violation databases from various states to help detect crime patterns will also add value to this crime detection paradigm.

Conclusion

The society we live in is a complicated and culturally revolutionised one, where crime problems are rising in an endless stream, and their prevention has become of utmost importance for the police and the government. In this project, we have applied data mining strategies, in particular, association rule mining strategies towards public security index requirement to support decision making for police departments and authorities.

42

Initially, an EDA was performed to demonstrate the baseline for understanding the Chicago crime dataset. Our EDA found many exciting results and statistics

that prompted us to apply rule mining algorithms on our dataset to uncover crime trends from the Chicago crime dataset. Before applying association rule mining, K Means clustering was carried out to understand the relationship between the three most important characteristics of the crime dataset, namely - location, time and crime type. Proper understanding of the graphs generated by clustering further warranted the appropriateness of association rule mining for this dataset. The most widely used association rule mining algorithm, Apriori, was first applied. Although the rules obtained using Apriori were reliable and had satisfactory values of lift/ interest and conviction, the application of Apriori was both time and computation intensive. To this end, a better association rule mining algorithm, FP Growth, was applied, which not only produced better results (rules with better conviction and lift) but also proved to be more time-efficient.

Finally, a GIS application, QGIS, was used to analyse the changing locations of hotspots with each week, month and year. Furthermore, a cluster& outlier analysis was conducted to identify low crime rate locations in high crime count areas to provide the police department with further insight into the reasons behind the crimes. ⁵In addition, a crime hotspot analysis was done to help in the deployment of police at most likely places of crime for any given window of time, to allow the most effective utilisation of police resources. Various results and graphs presented in this report corroborate the effectiveness of our model in identifying crime patterns.

43

Bibliography

Almanie, T., Mirza, R., and Lor, E. (2015). "Crime prediction based on crime types and using spatial and temporal criminal hotspots." arXiv preprint arXiv:1508.02050.

40 | Brown, D. E. and Hagen, S. (2003). "Data association methods with applications to law enforcement." Decision Support Systems, 34(4), 369–378.

Chandrasekar, A., Raj, A. S., and Kumar, P. (2015). "Crime prediction and classification in san francisco city." URL [http://cs229. stanford](http://cs229.stanford).^{7,8}

edu/proj2015/228 {_} report. pdf.

Dabbura, I. "K-means clustering: Algorithm, applications, evaluation methods, and draw- backs [online]. <https://towardsdatascience.com/k-means-clustering-algorithm-applications-evaluation-methods-and-drawbacks-aa03e644b48a>. Accessed On: Nov. 21 ,2019.

⁴¹ | esri. "How cluster and outlier analysis (anselin local moran's i) works [on- line]. <https://pro.arcgis.com/en/pro-app/tool-reference/spatial-statistics/h-how-cluster-and-outlier-analysis-anselin-local-m.htm>. Accessed On: Nov. 24 ,2019.

esri. "What is a z-score? what is a p-value? [online]. <http://desktop.arcgis.com/en/arcmap/10.3/tools/spatial-statistics-toolbox/what-is-a-z-score-what-is-a-p-value.htm>.

Iqbal, R., Murad, M. A. A., Mustapha, A., Panahy, P. H. S., and Khanahmadliravi, N. (2013). "An experimental study of classification algorithms for crime prediction." Indian Journal of Science and Technology, 6(3), 4219–4225.

Jain, V., Sharma, Y., Bhatia, A., and Arora, V. (2017). "Crime prediction using k-means⁹ algorithm." Global Research and Development journal¹⁰ for engineering Volume2, issue5.

Jangra, M. and Kalsi, M. S. (2019). "Crime analysis for multistate¹¹ network using naive bayes¹² classifier.

Keyvanpour, M. R., Javideh, M., and Ebrahimi, M. R. (2011). "Detecting and inves- tigating¹³ crime by means of¹⁴ data mining: a general crime matching framework." Procedia Computer Science, 3, 872–880.

⁴³ | Kianmehr, K. and Alhajj, R. (2008). "Effectiveness of support vector machine for crime hot-spots prediction." Applied Artificial Intelligence, 22(5), 433–458.

Kim, S., Joshi, P., Kalsi, P. S., and Taheri, P. (2018). "Crime analysis through machine learning." 2018 IEEE 9th Annual Information Technology, Electronics and Mobile Commu-¹⁵nication Conference (IEMCON)¹⁶, IEEE. 415–420.

44 | Li, S.-T., Kuo, S.-C., and Tsai, F.-C. (2010). "An intelligent decision-support model using ¹⁷fsom and rule extraction for crime prevention." Expert Systems with Applications, 37(10), 7108–7119.

44

45 | McClendon, L. and Meghanathan, N. (2015). "Using machine learning algorithms to ¹⁸an- alyze crime data." Machine Learning and Applications: An International Journal (MLAIJ), 2(1), 1–12.

46 | Nath, S. V. (2006). "Crime pattern detection using data mining." 2006 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology Workshops, IEEE. 41–44.

47 | Ng, V., Chan, S., Lau, D., and Ying, C. M. (2007). "Incremental mining for
48 | temporal association rules for crime pattern discoveries." Proceedings of the eighteenth conference on Australasian database-Volume 63, Australian Computer Society, Inc. 123–132.

Vaidya, O., Mitra, S., Kumbhar, R., Chavan, S., and Patil, R. (2018).

49 | "Comprehensive comparative analysis of methods for crime rate prediction.

50 | Yu, C.-H., Ward, M. W., Morabito, M., and Ding, W. (2011). "Crime forecasting using data mining techniques." 2011 IEEE 11th international conference on data mining ¹⁹work- shops, IEEE. 779–786.

45

1.	<i>was used</i>	Passive Voice Misuse	Clarity
2.	<i>was observed</i>	Passive Voice Misuse	Clarity
3.	<i>be aimed</i>	Passive Voice Misuse	Clarity
4.	<i>be increased</i>	Passive Voice Misuse	Clarity
5.	In addition → Also, Besides	Wordy Sentences	Clarity
6.	francisco → Francisco	Misspelled Words	Correctness
7.	stanford → Stanford	Improper Formatting	Correctness
8.	stanford → Stanford	Misspelled Words	Correctness
9.	the k-means, or a k-means	Determiner Use (a/an/the/this, etc.)	Correctness
10.	journal → Journal	Misspelled Words	Correctness
11.	a multistate, or the multistate	Determiner Use (a/an/the/this, etc.)	Correctness
12.	bayes → Bayes	Misspelled Words	Correctness
13.	inves-tigating → investigating	Confused Words	Correctness
14.	by means of → using, utilizing, employing, through	Wordy Sentences	Clarity
15.	Communi-cation → Communication	Confused Words	Correctness
16.	IEMCON → ICON	Misspelled Words	Correctness
17.	fsom → from	Misspelled Words	Correctness
18.	an-alyze → analyze	Confused Words	Correctness
19.	work-shops → workshops	Confused Words	Correctness

20.	<i>Deemed to be University under section 3 of UGC Act, 1956</i>	BITS Pilani Deemed to be University under Section 3 of UGC ... https://www.coursehero.com/file/p4rlmph/BITS-Pilani-Deemed-to-be-University-under-Section-3-of-UGC-Act-1956-Used-when/	Originality
21.	<i>We would like to express our deepest gratitude to our</i>	Household food insecurity access scale and dietary diversity score as a proxy indicator of nutritional status among people living with HIV/AIDS, Bahir Dar, Ethiopia, 2017	Originality
22.	<i>Many of these changes in the human condition have brought new challenges to the doorstep of the law</i>	The methods by which people and organizations choose to ... https://www.coursehero.com/file/p1806b6g/The-methods-by-which-people-and-organizations-choose-to-cope-with-change-will/	Originality
23.	<i>A framework of intelligent decision-support model based on a fuzzy</i>	An intelligent decision-support model using FSOM and rule ... https://www.sciencedirect.com/science/article/pii/S0957417410001855	Originality
24.	<i>attributed to a certain factor of randomness in the various crimes and the associated features; the branches of the decision trees are more</i>	Decision Stump algorithm could be attributed to a certain ... https://www.coursehero.com/file/p5um8q8/Decision-Stump-algorithm-could-be-attributed-to-a-certain-factor-of-randomness/	Originality
25.	<i>and give accurate results only if the test set follows the pattern modelled. On the other hand, the linear regression algorithm</i>	Figure 6 Selected Attribute Viewer Information Window ... https://www.coursehero.com/file/p4qef4m/Figure-6-Selected-Attribute-Viewer-Information-Window-Machine-Learning-and/	Originality
26.	<i>to take into account the dynamic nature of the</i>	Center of pressure (terrestrial locomotion) - Wikipedia	Originality

		https://en.wikipedia.org/wiki/Center_of_pressure_(terrestrial_locomotion)	
27.	Type ID Unique Identifier for the record Number Case Number The Chicago Police Department RD Number (Records Division Number), which is unique to the incident. Plain Text Date Date when the incident occurred. Date & Time Block The partially redacted address where the incident occurred, placing it o...	GitHub - k-chuang/chicago-crime-data-analysis: [Data ... https://github.com/k-chuang/Chicago-Crime-Data-Analysis	Originality
28.	Plain Text X Coordinate The x coordinate of the location where the incident occurred. Number Y Coordinate The y coordinate of the location where the incident occurred Number Year Year the Incident occurred Number Updated On Date and Time the record was last updated Date & Time Latitude The latitude...	GitHub - k-chuang/chicago-crime-data-analysis: [Data ... https://github.com/k-chuang/Chicago-Crime-Data-Analysis	Originality
29.	is shifted from the actual location for partial redaction but falls on the same block	GitHub - k-chuang/chicago-crime-data-analysis: [Data ... https://github.com/k-chuang/Chicago-Crime-Data-Analysis	Originality
30.	The algorithm works iteratively to assign each data point to one of the K groups based on the features that are provided. Data points are clustered based on feature similarity.	Biopython - Cluster Analysis - Tutorialspoint https://www.tutorialspoint.com/biopython/biopython_cluster_analysis.htm	Originality
31.	whose support count is greater than the minimum support count.	Association Rules - Saed Sayad https://www.saedsayad.com/association_rules.htm	Originality
32.	The COType field in the Output Feature Class will be	Optimized Outlier Analysis—Help ArcGIS Desktop http://desktop.arcgis.com/en/arcmap/latest/tools/spatial-statistics-toolbox/optimizedoutlieranalysis.htm	Originality

33.	<i>learn it from data, there is no right answer in terms of the number of clusters that we should have in any problem. Sometimes domain knowledge and intuition may help, but usually, that is not the case</i>	K-means Clustering: Algorithm, Applications, Evaluation ... https://towardsdatascience.com/k-means-clustering-algorithm-applications-evaluation-methods-and-drawbacks-aa03e644b48a	Originality
34.	<i>In the cluster-predict methodology, we can evaluate how well the models are performing based on different K clusters since clusters are used in the downstream</i>	K-means Clustering: Algorithm, Applications, Evaluation ... https://towardsdatascience.com/k-means-clustering-algorithm-applications-evaluation-methods-and-drawbacks-aa03e644b48a	Originality
35.	<i>The elbow method is one of the most</i>	Visitor Segmentation using K-means Clustering - Analytics ... https://medium.com/analytics-vidhya/visitor-segmentation-using-k-means-clustering-c874dcd41785	Originality
36.	<i>consistency within clusters of data [?]. The silhouette value is a measure of how similar an</i>	Silhouette Index – Cluster Validity index Set 2 ... https://www.geeksforgeeks.org/silhouette-index-cluster-validity-index-set-2/	Originality
37.	<i>The coefficient can take values in the interval [-1, 1]: If it is 0 → the sample is very close to the</i>	K-means Clustering: Algorithm, Applications, Evaluation ... https://towardsdatascience.com/k-means-clustering-algorithm-applications-evaluation-methods-and-drawbacks-aa03e644b48a	Originality
38.	<i>We have evaluated and compared the performance of</i>	Search results for: OCDMA https://publications.waset.org/abstracts/search?q=OCDMA	Originality
39.	<i>to see if there is a relationship between</i>	Answered: In a study to see if there is a... bartleby https://www.bartleby.com/questions-and-answers/in-a-study-to-see-if-there-is-a-relationship-	Originality

		between-students-consumption-of-alcoholic-beverages-and/6ce3baa6-b092-45e9-b379-a4060fb4c120	
40.	Hagen, S. (2003). "Data association methods with applications to law enforcement." <i>Decision Support Systems</i> , 34(4), 369–378.	Crime Data Mining, Threat Analysis and Prediction ... https://link.springer.com/chapter/10.1007/978-3-319-97181-0_9	Originality
41.	How cluster and outlier analysis (anselin local moran's i) works	Cluster and Outlier Analysis (Anselin Local Moran's I ... http://desktop.arcgis.com/en/arcmap/10.3/tools/spatial-statistics-toolbox/cluster-and-outlier-analysis-anselin-local-moran-s.htm	Originality
42.	An experimental study of classification algorithms for crime prediction." <i>Indian Journal of Science and Technology</i> , 6(3), 4219–4225.	Crime Data Mining, Threat Analysis and Prediction ... https://link.springer.com/chapter/10.1007/978-3-319-97181-0_9	Originality
43.	Alhaji, R. (2008). "Effectiveness of support vector machine for crime hot-spots prediction." <i>Applied Artificial Intelligence</i> , 22(5), 433–458.	Crime Data Mining, Threat Analysis and Prediction ... https://link.springer.com/chapter/10.1007/978-3-319-97181-0_9	Originality
44.	An intelligent decision-support model using fsm and rule extraction for crime prevention.	An intelligent decision-support model using FSOM and rule ... https://www.sciencedirect.com/science/article/pii/S0957417410001855	Originality
45.	crime data." <i>Machine Learning and Applications: An International Journal (MLAIJ)</i>	(PDF) USING MACHINE LEARNING ALGORITHMS TO ANALYZE CRIME ... https://www.academia.edu/12325398/USING_MACHINE_LEARNING_ALGORITHMS_TO_ANALYZE_CRIME_DATA	Originality
46.	2006 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology	Proceedings of the 2006 IEEE/WIC/ACM international ... https://dl.acm.org/citation.cfm?id=1194670	Originality

-
- | | | | |
|-------|--|--|-------------|
| 47. | <i>Ying, C. M. (2007). "Incremental mining for temporal association rules for crime pattern discoveries.</i> | Crime Data Mining, Threat Analysis and Prediction ...
https://link.springer.com/chapter/10.1007/978-3-319-97181-0_9 | Originality |
| <hr/> | | | |
| 48. | <i>Proceedings of the eighteenth conference on Australasian database-Volume</i> | Crime Data Mining, Threat Analysis and Prediction ...
https://link.springer.com/chapter/10.1007/978-3-319-97181-0_9 | Originality |
| <hr/> | | | |
| 49. | <i>Comprehensive comparative analysis of methods for crime rate prediction.</i> | (PDF) COMPREHENSIVE COMPARATIVE ANALYSIS OF METHODS FOR ...
https://www.academia.edu/36107249/COMPREHENSIVE_COMPARATIVE_ANALYSIS_OF_METHODS_FOR_CRIME_RATE_PREDICTION | Originality |
| <hr/> | | | |
| 50. | <i>Yu, C.-H., Ward, M. W., Morabito, M.</i> | Crime Data Mining, Threat Analysis and Prediction ...
https://link.springer.com/chapter/10.1007/978-3-319-97181-0_9 | Originality |
-