

Big Data Analytics and Mining for Effective Visualization and Trends Forecasting of Crime Data

Puppala Hemanth¹, Mareddy Vineeth Reddy²

IT Department, Guru Nanak Institutions Technical Campus, Hyderabad^{1,2}

Abstract: Enormous Data Analytics (BDA) is an efficient methodology for examining and recognizing various examples, relations and patterns inside a huge volume of information. In this paper we apply BDA to criminal information where exploratory information examination is directed for representation and patterns expectation. A few best-in-class information mining and profound learning methods are utilized. The prescient outcomes show that the Prophet model and Kera stateful LSTM perform in a way that is better than neural organization models, where the ideal size of the preparation information is discovered to be three years. These promising results will profit for police divisions and law requirement associations to more readily comprehend wrongdoing issues and give bits of knowledge that will empower them to follow exercises, foresee the probability of occurrences, adequately send assets and improve the dynamic cycle. With the help of such strategies, BDA can help us without any problem recognize wrongdoing designs which happen in a specific territory and how they are connected with time. The ramifications of AI and measurable methods on wrongdoing or other enormous information applications, for example, car crashes or time arrangement information, will empower the examination, extraction and comprehension of related examples and patterns, at last aiding wrongdoing anticipation and the executives. separating and standardization, Google maps-based Geo mapping of the highlights are actualized for perception of the measurable outcomes. Different methodologies in machine learning, profound learning, and time arrangement demonstrating are used for future patterns examination. 1) A progression of insightful investigations are directed to investigate and clarify the wrongdoing information in three US urban communities; 2) We propose a novel visual portrayal which is equipped for taking care of huge datasets and empowers clients to investigate, think about, and examine developmental patterns and examples of wrongdoing occurrences; 3) A mix and correlation of various AI, profound learning and time arrangement displaying calculations to foresee patterns with the ideal boundaries, time spans and models.

Keywords: Crime data Forecast , Visualization of Crime data, Big data analytics for crime data analytics.

I. INTRODUCTION

In most recent years, Big Data Analytics (BDA) has arisen as a rising method for reading records and extricating measurements and their individuals from the family in an extensive variety of programming regions. Due to non-stop urbanization and developing populations, towns play vital roles in our society. However, such traits have additionally been followed via way of means of a growth in violent crimes and accidents. To address such problems, sociologists, analysts, and protection establishments have dedicated lots attempt in the direction of mining ability styles and factors. In relation to public coverage however, there are numerous demanding situations in coping with massive quantities of to be had records. This gives, new approach plus technology wants to devised so as to research this heterogeneous and multi-sourced records. Analysis of such huge records allows us to efficiently hold tune of befall events, discover similarities from incidents, set up assets and make brief choices accordingly. This also can assist in addition our knowledge of each historic troubles and modern-day situations, in the long run making sure advanced protection/safety and first-class of life, in addition to expanded cultural and monetary increase.

The speedy increase of cloud computing and records acquisition and garage technology, from commercial enterprise and studies establishments to governments and numerous groups, have caused a massive quantity of exceptional scopes/complexities from records that has been accumulated and made publicly to be had. It has emerged as an increasing number of vital to extract significant statistics and attain new insights for knowledge styles from such records assets. BDA can efficiently cope with the demanding situations of records which can be too vast, too unstructured, and too speedy transferring to be controlled via way of means of conventional techniques. Like developing and influential system, BDA could resource groups to make useful of their records and facilitate advanced opportunities. Furthermore, BDA

may be deployed to assist sensible organizations circulate in advance with extra powerful operations, excessive earnings and glad customers. This can also help further our understanding of both historical issues and current situations, ultimately ensuring improved safety/security and quality of life, as well as increased cultural and economic growth. The rapid growth of cloud computing and data acquisition and storage technologies, from business and research institutions to governments and various organizations, have led to a huge number of unprecedented scopes/complexities from data that has been collected and made publicly available. It has become increasingly important to extract meaningful information and achieve new insights for understanding patterns from such data resources.

II. RELATED WORK

KUN NIU et al. has proposed. In this paper City-scale visitors velocity forecast gives considerable statistics basis to ITS, this advances suburbanites to the date statistics approximately visitors' condition. In any case, anticipating on-avenue car velocity correctly is trying, as the speed of vehicle on city avenue is tortured by numerous sorts of elements. These elements may be categorized into 3 fundamental perspectives, that are fleeting, spatial, and different inert insights. In this journal, we use Long Short-Term Memory. It is actually quite significant that our variant can avoid the extreme intricacy and vulnerability of emotional capacities extraction, and may be without difficulty prolonged to resolve different spatio-temporal forecast troubles together with waft expectation. The exploratory outcomes exhibit that the expectation variant we proposed can estimate city visitors' velocity effectively. J. Pera, A. Ferrández et al. The ODL stages have gotten outstandingly renowned lately. Conversations are a central particular mechanical assembly in various courses associated with online informative stages. These courses rely fundamentally upon discussion get-togethers for cooperation among students. In any case, the learning inclinations that these gadgets should give are constantly not abused. Conversations don't maintain learning assuming various messages are made, especially when they are posted in a befuddled and unstructured way which makes it irksome and drawn-out for the client to examine the data. Luka Stopar, Primo zSkriba et al. has proposed. presents a procedure for the instinctive discernment, examination and comprehension of huge multivariate time plan. Entrancing models with respect to such datasets by and large appear as periodic or dull direct much of the time achieved by the coordinated effort between factors. To recognize such models, we summarize the data as determined states, exhibiting common components as advances in the middle of state. The depiction could imagine huge information sets along possibly millions in models. In which, we loosen up depiction in various spatial granularities approving client in identifying plans on different scale. Yichuan Wang, Lee Ann Kung et al. has proposed. In this paper A significant data assessment engaged exchanged prototype ward.

Yuan-Yuan Liu et al. has proposed. In this paper the gauge of traveler numbers is huge for Destination Management and Marketing. While most existing procedures rely upon a lot of coordinated authentic data, using web search requests of the goal to assess its voyager appearances is one more way to deal with apply Big Data examination. In any case, there are no examinations exploring relationship of environment, temperatures, closures of the week and public events with the movement business objective appearances and web search inquiries of the goal, independently.

III. PROPOSED METHODOLOGY

Enormous Data Analytics (BDA) have transformed into appearing philosophy that breaks down the data and segregating information also the relationship to a larger degree in purpose zone. Equivalent to public game plan regardless, in which various troubles are there in dealing with tons of available information. Hence, fresh strategies and advancements should be made to review this heterogeneous and multi-got information. Subsequently, new procedures and developments ought to be created to inspect this heterogeneous and multi obtained data. The start of tremendous information in BDA, the assessment and connected troubles while imparting them. On investigation, openings and challenges of bad behavior in data mining. In extra to that, this endeavor information about the data searching for observing the model and examples in bad behavior to be used fittingly and to be a help for fledglings in the assessment of bad behavior data mining. As a result, the organization and the examination with huge data are particularly irksome and complex. Constructing adequacy in bad behavior disclosure, it is critical to pick the data mining methods sensibly. different data mining application. Besides, couple of techniques that has made separate connection between two item sets even more enough, for instance, normal information thought anyway the computation was extended the more proportion of time.

IV. CRIME DATA MINING, VISUALIZATION & TRENDS FORECASTING

In criminology literature the relationship between crime and various factors has been intensively analyzed, where typical examples include historical crime records [4], unemployment rate [5], and spatial similarity [6]. Using data mining and statistical techniques, new algorithms and systems have been developed along with new types of data. For instance, classification and statistical models are applied for mining of crime patterns and crime prediction [8], where transfer

learning has been employed to exploit spatio-temporal patterns in New York city [9]. Wu et al. [2] developed a system to automatically collect crime-logged data for mining of crime patterns, in order to achieve more effective crime prevention in/around a university campus. In Vineeth et al. [1], a random forest was applied on the obtained correlation between crime types to classify the state based on their crime intensity point. Unsupervised learning based methods have also been used for mining of crime patterns and crime hotspots, such as memetic differential fuzzy cluster for forecasting of criminal patterns, and fuzzy C-means algorithm to cluster criminal events in space. Noor et al. [4] derived association mining rules to determine relationships between different crimes.

Injadat et al. [5] conducted a survey to summary data mining techniques on social media. With deep learning and neural networks, new models have been developed to predict crime occurrence [25]. As deep learning and artificial intelligence have achieved great success in computer vision, they are also applied in BDA for predicting trends and classification. In Zhao et al. [3] and Dai et al. [3], Long Short-Term Memory networks were successfully applied to predict stock price and gas dissolved in power transformation. In Kashef et al. [3], a neural network was applied on a smart grid system to estimate the trends of power loss. Zheng et al. [5] proposed a big data processing architecture for radio signals analysis. Zhao et al. proposed a neural network model to predict travel time and gained high accuracy. Niu et al. utilized LSTM and developed an effective speed prediction model to solve spatio-temporal prediction problems. Peral et al. summarized analytics techniques and proposed an architecture for online forum data mining.

In our proposed system, a similar but more comprehensive workflow has been adopted, which include statistical analysis, data visualization, and trends prediction. Data visualization and mining techniques are used to show the extracted statistical relationships among different attributes within the huge volume of data. State-of-the-art machine learning and deep learning algorithms are deployed to forecast trends and obtain optimal models with the highest accuracy.

V. DATA ANALYSIS AND VISUALIZATION

The three crime datasets we used for analysis are publicly available, which cover 3 cities in US, i.e. San-Francisco, Chicago, and Philadelphia. The San-Francisco crime data contains 2,142,685 crime incidents from 01/01/2003 to 11/08/2017 [38]. Data from Chicago has a total number of 5,541,398 records, dating back from 2017 to 2003 [39]. In the Philadelphia dataset, there are 2,371,416 crime incidents which were captured from 01/01/2006 to 12/31/2017 [40].

Detailed analysis of these dataset is presented as follows.

FEATURED ATTRIBUTES

For each entry of crime incidents in the datasets, the following 13 featured attributes are included:

- 1) IncidentNum - Case number of each incident;
- 2) Dates - Date and timestamp of the crime incident;
- 3) Category - Type of the crime. This is the target/label that we need to predict in the classification stage;
- 4) Descript - A brief note describing any pertinent details of the crime;
- 5) DayOfWeek - Day of the week that crime occurred;
- 6) PdDistrict - Police Department District ID where the crime is assigned;
- 7) Resolution - How the crime incident was resolved (with the perpetrator being, say, arrest or booked);
- 8) Address - The approximate street address of the crime incident;
- 9) X - Longitude of the location of a crime;
- 10) Y - Latitude of the location of a crime;
- 11) Coordinate - Pairs of Longitude and Latitude;
- 12) Dome - whether crime id domestic or not;
- 13) Arrest - Arrested or not;

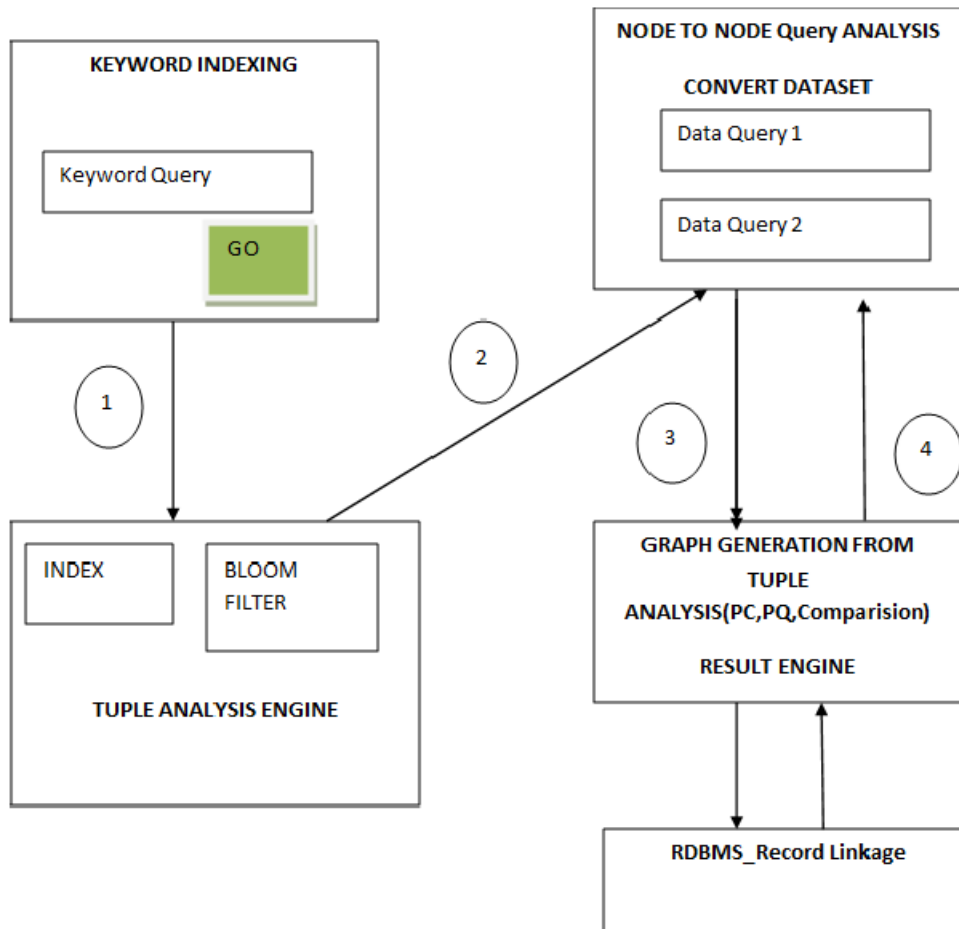


Fig 1. Architecture Diagram

Information perception (oftentimes contracted information viz) is an interdisciplinary discipline that proposals with the picture delineation of data. It is a especially green manner of speaking while the data is cut off as an illustration a Time Series. From a scholastic component of view, this illustration may be taken into thought as a planning among the authentic information (by and large mathematical) and picture components Information perception has its foundations withinside the discipline of Statistics and is consequently normally taken into thought a branch of Descriptive Statistics. Notwithstanding, because of the fact each layout abilities and factual and figuring abilities are expected to visualize efficiently, its miles contended through method of means of a few creators that it's miles each an Art and a Science. To talk statistics genuinely and effectively, data perception utilizes measurable photographs, plots, statistics photos and various devices. Mathematical information can be encoded the utilization of specks, follows, or bars, to talk a quantitative message outwardly.

Compelling representation allows customers dissect and motive approximately information and proof. It makes complicated information extra available, fathomable and usable. Clients could likewise furthermore have explicit scientific undertakings, like making examinations or information causality, and the layout precept of the picture (i.e., showing correlations or showing causality) follows the errand. Tables are typically utilized where in customers will appearance up a chose estimation, simultaneously as outlines of numerous sorts are utilized to reveal styles or connections withinside the information for one or additional factors. Information representation alludes back to the methods used to talk information or statistics via method of method for encoding it as noticeable items (e.g., factors, follows or bars) contained in photographs. The expectation is to talk statistics genuinely and effectively to clients. It is one of the means in information evaluation or data science. As per Vitaly Friedman (2008) the "important intention of data representation is to talk statistics genuinely and efficiently via graphical means. It would not imply that data representation wishes to appearance uninteresting to be useful or extraordinarily state-of-the-workmanship to appearance lovely. To deliver thoughts efficiently, every stylish shape and usefulness want to move connected at the hip, granting bits of knowledge directly into a as a substitute inadequate and complicated information set through method of means of speaking its key-factors in an extra natural way. However, planners much of the time neglect to acquire stability among shape and work,

making appropriate information perceptions which neglect to fill their significant need - to talk statistics". For sure, Fernanda Viegas and Martin M. Wattenberg forewarned that an incredible perception need to know no longer simple stalk genuinely, yet invigorate watcher commitment and consideration. Information perception is firmly related with statistics photos, insights representation, clinical representation, exploratory information evaluation and measurable photographs. In the new thousand years, data representation has end up an energetic place of examination, instructing and improvement. As per Post et al. (2002), it has joined clinical and measurements perception. In the industrial surroundings information visualization is frequently called dashboards. Data photographs are each other exceptionally typical place shape of data visualization.

VI. PREDICTION MODELS

In order to tackle the problem of crime trends forecasting we explored several state-of-the-art machine learning and deep learning algorithms and time series models. A time series is a sequence of numerical data points successively indexed or listed/graphed in the time order. Usually, the successive data points within a time series are equally spaced in time, hence these data are discrete in time. Fig. 8 demonstrates how the amount of crime incidents changed over time, which clearly show the potential trend and seasonality in the data as analyzed and discussed in the following sections.

PROPHET MODEL

The Prophet model is a procedure for forecasting time series data based on an additive model where non-linear trends are t with yearly, weekly, and/or daily seasonality, plus holiday effects [36]. It works best with time series that have strong seasonal effects and cover several seasons of historical data. The Prophet model is robust to missing data and shifts in the trend, and typically it handles outliers well. The Prophet model is designed to handle complex features in time series, it also designed to have intuitive parameters that can be adjusted without knowing the details of the underlying model. Time series can exhibit a variety of patterns, and it is always helpful to decompose a time series into several components, each representing an underlying pattern category. Fig. 8 illustrates a decomposed crime time series, where for each original time series on the top, the three decomposed parts can respectively show the estimated trend component, seasonal component, and irregular component, respectively. The estimated trend component has shown that the overall crimes in San-Francisco slightly decreased from 2003 to 2013, followed by a steady increase from then on to 2017. However, crimes in Chicago seemed to decrease quickly from 2003 to 2015 and then became quite stable, whilst in Philadelphia the number of crimes had a downward trend yet with some undulations until it became stable after 2016. Regarding the seasonal component, it changes slowly over the time, where a quite strong annual periodic pattern can be observed for Chicago and Philadelphia than San-Francisco. This may be due to the more apparent annual climate changes in the two cities, where the crimes reached the peak in the middle of the year when the temperature becomes the hottest. As such, time series models will perform well on these datasets to forecast crimes in the future.

Years for training	RMSE-Prophet	Correlation-Prophet	RMSE-LSTM	Correlation-LSTM	RMSE-Neural Network	Correlation-Neural Network
10	38.21	0.384	55.21	0.354	54.79	0.097
5	35.70	0.402	54.18	0.365	48.04	0.232
4	36.18	0.415	53.96	0.411	48.04	0.236
3	35.65	0.398	45.65	0.423	41.62	0.291
2	91.93	0.087	160.2	0.098	41.17	0.128
1	100.56	0.182	95.96	0.122	-	-
10	76.89	0.560	77.01	0.532	77.19	0.367
5	68.21	0.652	69.12	0.549	92.74	0.551
4	66.75	0.654	67.45	0.612	88.04	0.492
3	66.68	0.658	67.15	0.625	75.06	0.505
2	67.42	0.632	68.14	0.576	75.90	0.552
1	100.51	0.02	78.98	0.459	-	-
10	51.83	0.716	50.65	0.709	82.21	0.422
5	56.07	0.728	51.23	0.698	71.19	0.486
4	55.37	0.728	49.22	0.714	67.74	0.588
3	48.73	0.729	48.15	0.725	63.68	0.537
2	50.35	0.718	57.16	0.705	170.12	0.128
1	100.71	0.098	140.63	0.562	-	-

Table 1. Comparison of different algorithms/models in terms of RMSE and spearman correlation under different sizes of training samples.

To train our models for predicting trends, we first summarized the number of crime incidents per day, and then transformed these data into a "tibbletime" format, and then we divided the data into training and testing sets, where the training set contains data from 2003 to 2016 and the testing set has data from 2017, for training process we set 1 year's data as validation set. We evaluated the performance of the prediction models whilst changing the number of training years from 1 to 10 and the results are summarized in Table 1. As seen in Table 2, more training data do not necessarily lead to better results although too little training data also fails to generate good results. The optimal time period for crime trends forecasting is 3 years where the RMSE is the minimum and the spearman correlation is the highest.

VII. LSTM MODEL

LSTM model is a powerful type of recurrent neural network (RNN), capable of learning long-term dependencies. For time series involves auto-correlation, i.e. the presence of correlation between the time series and lagged versions of itself, LSTMs are particular useful in prediction due to their capability of maintaining the state whilst recognizing patterns over the time series. The recurrent architecture enables the states to be persisted, or communicate between updated weights as each epoch progresses. Moreover, the LSTM cell architecture can enhance the RNN by enabling long term persistence in addition to short term.

$$\begin{aligned}f(t) &= \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \\i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\ \tilde{C}_t &= \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \\ C_t &= f_t * C_{t-1} + i_t * \tilde{C}_t\end{aligned}$$

where, f_t is a sigmoid function to indicate whether to keep the previous state, C_{t-1} is the old cell state, C_t is the updated cell state, W_f , W_i , and W_C are the previous value in each layer, h_{t-1} and x_t is the input value, b_f , b_i , and b_C are constant values, it decides which value will be used to update the state, C_t stands for the new candidate values.

VIII. CONCLUSION & FUTURE WORK

In this paper a series of state-of-the-art big data analytics and visualization techniques were utilized to analyze crime big data from three US cities, which allowed us to identify patterns and obtain trends. By exploring the Prophet model, a neural network model, and the deep learning algorithm LSTM, we found that both the Prophet model and the LSTM algorithm perform better than conventional neural network models. We also found the optimal time period for the training sample to be 3 years, in order to achieve the best prediction of trends in terms of RMSE and spearman correlation. Optimal parameters for the Prophet and the LSTM models are also determined. Additional results explained earlier will provide new insights into crime trends and will assist both police departments and law enforcement agencies in their decision making. In future, we plan to complete our on-going platform for generic big data analytics which will be capable of processing various types of data for a wide range of applications. We also plan to incorporate multivariate visualization graph mining techniques and fine grained spatial analysis to uncover more potential patterns and trends within these datasets. Moreover, we aim to conduct more realistic case studies to further evaluate the effectiveness and scalability of the different models in our system.

REFERENCES

- [1]. Niu K, Zhang H, Zhou T, et al. A Novel Spatio-Temporal Model for City-Scale Traffic Speed Prediction. IEEE Access, vol. 7, pp. 30050 - 30057, Feb. 2019.
- [2]. Peral J, Ferrández A, Mora H, et al. A Review of the Analytics Techniques for an Efficient Management of Online Forums: an Architecture Proposal. IEEE Access, vol. 7, pp. 12220 - 12240, Feb. 2019.
- [3]. Stopar L, Skraba P, Grobelnik M, et al. Streamstory: exploring multivariate time series on multiple scales. IEEE transactions on visualization and computer graphics, vol. 25, no. 4, pp. 1788-1802, Apr. 2019.
- [4] Ravindra Changala, "EVALUATION AND ANALYSIS OF DISCOVERED PATTERNS USING PATTERN CLASSIFICATION METHODS IN TEXT MINING" in ARPN Journal of Engineering and Applied Sciences, Volume 13, Issue 11, Pages 3706-3717 with ISSN:1819-6608 in June 2018.

- [5] Y.Wang, L. Kung,W. Y. C.Wang, and C. G. Cegielski, "An integrated big data analytics-enabled transformation model: Application to health care," *Inf. Manage.*, vol. 55, no. 1, pp. 6479, Jan. 2018.
- [6] J. Archenaa and E. A. M. Anita, "A survey of big data analytics in healthcare and government," *Procedia Comput. Sci.*, vol. 50, pp. 408-413, Apr. 2015.
- [7] Ravindra Changala, "Pattern Deploying Methods for Text Mining" in *International Journal of Soft Computing*, Volume 13, Issue 2, pages 61-68 with ISSN: 1816-9503 in June 2018.
- [8] W. Grady, H. Parker, and A. Payne, "Agile big data analytics: Analytics Ops for data science," in *Proc. IEEE Int. Conf. Big Data*, Boston, MA, USA, Dec. 2017, pp. 2331-2339.
- [9] Y.-Y. Liu, F.-M. Tseng, and Y.-H. Tseng, "Big Data analytics for forecasting tourism destination arrivals with the applied Vector Autoregression model," *Technol. Forecasting Social Change*, vol. 130, pp. 123-134, May 2018.
- [10] S. Musa, "Smart cities: A road map for development," *IEEE Potentials*, vol. 37, no. 2, pp. 19-23, Mar./Apr. 2018.
- [11] Ravindra Changala, "Classification by Decision Tree Induction Algorithm to Learn Decision Trees from the class-Labeled Training Tuples" published in *International Journal of Advanced Research in Computer Science and Software Engineering (IJARCSSE)*, ISSN: 2277 128X , Volume 2, Issue 4, April 2012.
- [12] Ravindra Changala, "Decision Tree Induction Approach for Data Classification Using Peano Count Trees" published in *International Journal of Advanced Research in Computer Science and Software Engineering (IJARCSSE)*, ISSN: 2277 128X, Volume 2, Issue 4, April 2012.
- [13] S.Wang, X.Wang, P. Ye, Y. Yuan, S. Liu, and F.-Y.Wang, "Parallel crime scene analysis based on ACP approach," *IEEE Trans. Computat. Social Syst.*, vol. 5, no. 1, pp. 244255, Mar. 2018.
- [14] S. Yadav, A. Yadav, R. Vishwakarma, N. Yadav, and M. Timbadia, "Crime pattern detection, analysis & prediction," in *Proc. IEEE Int. Conf. Electron., Commun. Aersp. Technol.*, Coimbatore, India, Apr. 2017, pp. 225-230.
- [15] Wang Z, Ren J, et al.: A deep-learning based feature hybrid framework for spatiotemporal saliency detection inside videos. *Neurocomputing*, vol. 287, pp. 68-83, Apr. 2018.
- [16]. Yan Y, Ren J, et al.: Cognitive Fusion of Thermal and Visible Imagery for Effective Detection and Tracking of Pedestrians in Videos. *Cognitive Computation*, vol. 10, no. 1, pp. 94-104, Feb. 2018.
- [17] Ravindra Changala, "A SURVEY ON DEVELOPMENT OF PATTERN EVOLVING MODEL FOR DISCOVERY OF PATTERNS IN TEXT MINING USING DATA MINING TECHNIQUES" in *Journal of Theoretical and Applied Information Technology* in 31st August 2017. Vol.95. No.16, ISSN: 1817-3195, pp.3974-3987.
- [18]. Dai J., Song H. et al., LSTM networks for the trend prediction of gases dissolved in power transformer insulation oil, 2018 12th Int. Conf. on the Properties and Applications of Dielectric Materials , Xi'an, China, 2018, pp. 666-669.