# An Overview of Data Mining Techniques and its Realtime Applications

## K. Chitra Lekha[1*], S. Prakasam[2]

[1]Department of Computer Science and Applications, SCSVMV University, Kanchipuram, India
[2]Department of Computer Science and Applications, SCSVMV University, Kanchipuram, India

*Corresponding Author:   chithuparthi03@gmail.com

*Abstract*— Data mining, also prevalently referred as Knowledge innovation from data, is the robotic or expedient withdrawal of patterns representing information implicitly hoard or confined in huge databases, data warehouses, the web, the data streams or other enormous information repositories. Data mining is the expertise that congregates up to the dispute of solving our pursuit for acquaintance from these cosmic data burdens. It affords us with a client oriented loom to novel concealed prototypes in data. This paper accomplishes a prescribed assessment of the perception of data-mining, the typical tasks engross in data-mining, its relevance in day to day field, techniques and methodology. In additional to that, this paper affords insight for concerning the data mining to vindicate the patterns and trends to be utilized suitably and to be a supportive for beginners in the research of  data mining. The core intention of this manuscript is to congregate more interior perceptions and skills in data mining.

*Keywords*—Artificial Intelligence, Time series analysis, Regression, Prediction

## I. INTRODUCTION

Data mining proffers the programmed, eventual investigation beyond the analyses of precedent events provided by retrospective tools emblematic of decision support systems. Data mining is one among the in numerous amount of investigative tools for evaluating data. It permits the users to evaluate data from several diverse scopes or angles, sort it, and recapitulate the relationships recognized. Precisely, data mining is the progression of ruling correlations or patterns amid dozens of fields in huge relational databases. It is intricate to acquire quicker and suitable information by conventional physical investigation which is monotonous as well as very burdensome. As a result, data mining is utilized on the whole i) to shrink costs in the course of appropriate recognition and preclusion of ravage and scam, ii) acquiring suitable and up-to-date information and iii) enhance revenues through better advertising approach. Despite the fact that data mining techniques are a means to constrain efficiencies and envisage client deeds, if utilized acceptably, a trade can put itself, apart from its rivalry through the exploit of prognostic investigation. Explicit data mining benefits vary depending on the aspiration and the industry [1].

Data mining or data discovery is that the system aided technique of excavating through and investigating immense sets of information so extorting which means of the tidings. Data dispensation tools envisage behaviours and prospect trends, consenting businesses to shape practical, knowledge-driven assortments [2].

**KDD:** Knowledge discovery in databases is a rapidly emergent domain, whose progress is obsessed by well-built explore interests as well as insistent, convenient, communal, and economical desires. While in the proceeding few years, knowledge detection tools have been utilized essentially in research environments and classy software products are currently emerging quickly. Knowledge innovation in databases is the nontrivial progression of recognizing legal, novel, potentially valuable, and eventually explicable prototypes in data [3]

**a. Pre-mining tasks:**
**Data cleaning:** This is the initial pace to eradicate noise data and extraneous data from collected untreated data.

**Data integration:** At this pace, diverse data resources are united into significant and practical data so that it may be utilized for further handling of data.

**Data transformation:** In this stride, by performing diverse processes such as smoothing, normalization or summarization, the data is transformed or united into mandatory forms for mining.

**b. Post-mining tasks:**
**Pattern evaluation:**
At this stride, based on prearranged trials, smart prototypes representing facts and information are recognized.

**Knowledge representation:**
This is the final phase in which, revelation and acquaintance depiction skills are utilized to facilitate clients to recognize and deduce the data mining facts or consequences. The ultimate target of knowledge innovation and data mining progression is to discern the prototype that are anonymous among the enormous set of data and deduce valuable acquaintance and information. The motivation for proceeding with this review work is to abet a serving hand over to the young researchers who are undergoing their study in the domain of data mining.

The remaining of the manuscript is organized as follows Section II describes the Roots of data mining, Section III contain the Needs and Goals of data mining, Section IV illustrates the Real time applications of data mining, Section V describes the Disputes of data mining, Section VI contains merits of data mining whereas Section VII contain the demerits of data mining and Section VIII concludes research work with future directions.

## II. ROOTS OF DATA MINING

Data Mining is the course of action which embraces appraising and probing hefty pre-existing databases in order to spawn novel information which might be vital to the society. The extraction of innovative information is envisaged by utilizing the existing datasets. Numerous approaches for investigations and prophecy in data mining had been executed [4]

### 2.1 Statistics
Devoid of statistics, there would be no data mining, since statistics are the basis of nearly all skills on which data mining is fabricated. Statistics grip conceptions such as regression scrutiny, normal distribution, standard deviation, standard variance, discriminate examination, cluster investigation, along with confidence intervals, all of which are utilized to cram fact and data associations.

### 2.2 Artificial Intelligence & Machine Learning
Data mining's subsequent best ever kin line is artificial intelligence and machine learning. AI is built ahead heuristics as conflicting to statistics, and endeavours to pertain human-thought like dispensation to statistical tribulations. Machine Learning may perhaps be measured as a progression of AI, since it mingles AI heuristics with enhanced statistical techniques. Let system programs be trained to gain knowledge of concerning the data they cram and then pertain learned knowledge to data.

### 2.3 Databases
Enormous quantity of information needs to be hoarded in a depository, and that too desires to be administered. Data warehousing also supports OLAP manoeuvres to be functional on it, to prop up decision making. The progress in

Information Technology has spawned huge quantity of databases and massive data in diverse domains. The exploration in databases and information technology has given augment to an approach to hoard and operate this valuable data for auxiliary decision making [5].

### 2.4 Datamining Algorithms & Techniques
**2.4.1 Predictive:** The **predictive form** executes by edifying a prophecy about ideals of data, which utilizes notorious consequences initiated from diverse datasets.

**2.4.1.1 Classification:** Classification is the most generally practical data mining method, which utilizes a set of pre-classified paradigms to build up a replica that can catalogue the population of reports at hefty. This approach regularly utilizes decision tree or neural network-based classifier algorithms. The widespread distinctiveness of classifier tasks are as supervised learning, sorts reliant variable and conveying innovative data to one of a set of well-defined modules. E.g. Given classes of patients that are in contact to medical treatment retorts, discover the type of healing to which an innovative patient is most liable to counter.

**2.4.1.2 Regression:**
Regression is an additional prognostic data-mining model is also recognized as supervised learning system. This system evaluates the reliance of several attribute values, which is reliant ahead the values of further attributes generally there in similar item. The intention significance values are known in the regression models. E.g. based on ancestors' records, the deeds of child might be envisaged by utilizing a regression paradigm.

**2.4.1.3 Timeseries Analysis:**
Time-series database utilizes progression of values or proceedings attained over repetitive dimensions of time. The ideals are normally summarized at identical intervals of time such as hourly, weekly, daily. A sequence database is one database that embraces progression of prearranged trials, occasionally having tangible designs of time. Estimating certain crucial values over enduring years might be probable by means of the time series data analysis. This can escort to enhanced evaluation of potential chucks and also scheduling for prospective enlargement. E.g. Stock market.

**2.4.1.4 Prediction:**
This method discerns the correlation involving independent variables and the association amid reliant and independent variables. The prediction is to envisage a prospect state, rather than a recent solitary. E.g. Given a prophetic replica of credit card dealings, envisage the probability that a precise transaction is fake or legitimate. Prediction can also be utilized to authenticate a revealed hypothesis.

**2.4.2 Descriptive:** A **Descriptive form** presents the information in brief outline which is effectively a synopsis of the data points, uncovers prototypes in the data and recognizes the interactions involving attributes signified by the data.

### 2.4.2.1Clustering:

It is an unsupervised model innovative approach where the data are assembled together based on a resemblance appraisal. It is way of vindicating resemblance between data according to their traits. Clustering can be considered as recognition of analogous modules of objects. By utilizing clustering techniques, we can further spot intense and bare areas in object space and can discern overall allocation model and associations among data characteristics.
E.g. Given a data set of consumers discovers subgroups of consumers who have related trade manners.

### 2.4.1.2Sequence Discovery:

It exposes relationship amid data. It is set of entities each connected with its own timeline of events. E.g. scientific trials, expected catastrophe like floods and investigation of DNA progression.

### 2.4.1.3Summarization:

Summarization is characterized as the abstraction or generality of statistics. The summarization system plots facts into subsets with plain metaphors. The recapitulated data set confers broad outline of the data with cumulative information. Simple summarization techniques such as tabularizing the significance and standard deviations are frequently applied for data scrutiny, data revelation and computerized report generation. E.g. The total minutes, seconds and height covered by athletes in protracted distance race can be summarized.

### 2.4.1.4Association rules:

The Association technique is utilized to extort the interactions involving elements and objects. In this method, the occurrence of solitary replica implies the existence of another form i.e. an item is correlated to another in stipulations of cause-and-effect. This is general in ascertaining a type of arithmetical interactions between diverse co-dependent variables of data mining; association rules are handy for exploring and forecasting patron deeds. They moreover play a significant task in shopping basket data analysis, product clustering and directory design and stockpile outline.

### 2.5Visualization Techniques

Visualization techniques are awfully valuable routine of discerning prototypes in data sets, and might be utilized at the commencement of a data mining progression. There is an intact domain of exploration enthusiastic to investigate for fascinating protrusions of datasets – this is termed as

Projection Pursuit. For instance, clusters are frequently arithmetical signified. Moreover, a huge set of conventions is simpler to recognize when controlled in a hierarchical trend and graphically observed such as in the structure of a decision tree. Visualization methods may perhaps vary from plain scatter plots and histogram designs over parallel coordinates to 3D movies.

### III.    NEEDS & GOALS OF DATA MINING

#### 3.1 NEEDS OF DATA MINING:

The accomplishment of digital revolt and the appreciation of the internet have brought about a enormous quantity of multi-dimensional statistics in practically all human venture, and the data type varies from images, text, audio, graphics, speech, video and hypertext thus providing organizations with too many data, but the whole data might not be valuable if it does not afford a substantial distinctive information that might be utilized in unravelling a quandary [6]. The pursuit to spawn information from accessible data provoked necessitate for data mining.

#### 3.2 Goals of Data Mining:

Data mining is principally done with aspire of attaining firm objectives and it range from clustering, association, optimization, classification, prediction and identification. The goals of data mining are fast retrieval of discovery of acquaintance from the databases, information or data to recognize concealed patterns and those prototypes which are formerly not investigated, to shrink the intensity of intricacy, reduction of time etc [7].

### IV.    REALTIME APPLICATIONS OF DATA MINING

An assortment of fields personalized data mining expertise because of rapid access of data and precious information from a hefty amount of facts.

**Supply chain visibility:**
Companies have robotic segments of their supply chain facilitating compilation of momentous data about inventory, supply concert and logistic of equipments, and refined wares, material expenditures, precision of tactics for order liberation. Data mining relevance also extends through cost optimization and efforts oblige scrutiny in institutes.

**Geospatial decision making:**
In case of typical weather facts and soil ecosystem scenario, habitual extraction and analysis of fascinating patterns concerning modelling environmental data and conniving, a proficient algorithm for vindicating spatiotemporal prototypes in the form of telex-connection patterns or recurring and persistent climatic patterns is in need. This manoeuvre is usually conceded out utilizing the clustering technique, which splits the data into significant factions and

facilitates to computerize the innovation of telex-connections.

**Biomedicine and science application:**
Biology with the crash of data mining has progressed into a meadow of immense science attitude concerning and accumulating data, extract for novel hypothesis then validate with data or supplemental research. It as well embraces innovation of prototypes in radiological metaphors, scrutiny of microarray (gene-chip) experimental facts to cluster genes and to transmit to indication or infection, investigation of side effects of drugs and efficacy of certain drugs.

**Forensic & scandalous exploration:** Data mining skill is utilized in forensic and enforcement sector to levy prior scandalous reports in order to recognize the criminal as well as to conclude the crime outline ,attitudes of the culprits and the indict.

**Data mining in banking and investment:** Data mining has been utilized widely in the depository and fiscal markets. In the banking domain, data mining is utilized to envisage credit card scam, to guess peril, to evaluate the trend and productivity. In the financial markets, data mining method such as neural networks are utilized in stock forecasting, penalty prophecy and so on. Data mining proposals are exploited in many banking domains for consumer segmentation and production, tribute attains and consents, envisaging expense default, marketing, perceiving sham transactions [8]

**Data Mining in Telecommunication**: Owing to huge availability of data, rapidly altering and exceedingly aggressive situation, telecom networks utilize data mining techniques to progress their advertising efforts, recognition of scam, and enhanced administration of telecommunication network [9].

**Data Mining in Cloud Computing:** The utilization of data mining techniques through Cloud computing will consent the clients to recover significant information from practically incorporated data warehouse that shrinks the expenses of infrastructure and storage [10].

**Data mining in detection of Cyber crimes:** Cyber crimes are a mutual pest and charge our humanity greatly in several customs. The exploration of cyber crime cases have extremely considerable task in law enforcement system in any nation. Data mining skill is applied to con detection to determine the sting detection model, to portray the progression of creating the scam detection model, and then to instigate the data model with any classifier [11]. Exploitation of data mining skills in the revealing of cyber crime can amend the circumstances of decision makers and law enforcement officials in an enhanced mode. Data mining

plays a decisive role for captivating assessments on numerous concerns correlated to prophecy of cyber crime [12].

## V. DISPUTES OF DATA MINING

Data mining can be utilized to discover prototypes and correlations that would otherwise be thorny to stumble on. This skill is trendy with several dealings since it permits them to discover more about their clients and formulate elegant marketing assessments [13]. The foremost disputes in data mining are: Data mining has to transact with massive quantity of facts situated at diverse spots; the quantity of data can straightforwardly go beyond the terabyte edge; data mining is concerts [14]. The concerns in data mining are emerging day by day, and researchers are frustrating to progress and formulate data competent by applying appropriate algorithm. In future, data mining should recognize the intricate inputs from the clients and must spawn the valuable and required outputs [15]. very computationally rigorous process concerning very huge data sets. Generally, it is obligatory to panel and allocate the data for parallel dispensation to attain suitable time and space

## VI. MERITS OF DATA MINING

- It envisages prospect drifts.
- It facilitates in production of assessments.
- It assists to progress the proceeds of the private companies.
- It is primarily utilized in market scrutiny.
- It is efficiently exploited in detection of scam.
- DM skills are handy for health care insurers to detect hoax and neglect.
- DM facilitates the surgeons to discover effectual treatments and preeminent observations through healthcare software.
- In case of huge datasets, it is probable to fasten the task by utilizing data mining functionalities.
- DM aids in production of faster reports and more rapidly investigation, which force increase in equipped competence and also shrinks operating expenditure.

## VII. DEMERITS OF DATA MINING

- The intricacy and heterogeneity of in the volume of data would generate redundant statistical cataloging.
- Have to consider ethical, permissible along with societal concerns.
- The law should suit while trading with data possession.
- In case of medical data, issues might arise in privacy, security of human data administrative concerns.
- Maltreatment of information or inaccurate information.

**Table1. Data mining methods, Tools, Applications and Constraints**

| Data mining method | Explanation | Implementation / Tools | Applications | Constraints |
|---|---|---|---|---|
| Classification | Consortium of data items, Generation of factions with similar aspects and distinctiveness | Orange, WEKA, SPSS, Python, R etc | Diminution of information intricacy, Reformation in data compilation, courteous in forecasting. | Ineffectual for diverse data |
| Statistics | Tabular and graphical depiction of information | Analog, SPSS, WebStat, AccessWatch, | Evaluating the number of visits on the network, Huge information is depicted as tables, grids, 3D representation. | Fail to scrutinize individual points, Not executed on diverse data, A minute fault might lead to ambiguous. |
| Clustering | Federation of information by resemblance | PYCLUSTER, Cluster 3.0, Java Tree View etc. | Liability lenience, preservation issues | Doesn't prop up mutual storage, equipped howlers |
| Web data mining | Data mining associated to network | CrawalMonster, Winautomation, import.io etc. | Deceiving the web configuration, contents of the web and convention of networks. | Incursion of solitudes, Extraneous stuffing. |
| Correlation | Conception of prototypes from gestures, audios, cartridges, images, sequences | Google Correlate, Google Trends, Google Flu Trends | Very helpful for researchers to accumulate new data than experiments, Utilized in neuroscience, fabric skill and economics | Doesn't afford the motive of relation between items |
| Outlier detection | Scrutiny of one data item that is unnecessary from a group of other data items | CMSR Data Miner | Dispensation of images, revealing of manufacturing dents, scam exposure, invasion revealing, recognition of indoor trading, therapeutic and public physical condition | Require probabilistic data model that Is intricate and also necessitate statistical validation. |
| Association Rule Mining | Antecedent as well as consequent statements | KNIME, SPSS, Magnum Opus, WEKA, FrIDA, FPM, ORANGE, Bart Goethals | Dealing of stockpile, Exploited in Leave administration system. | Algorithm utilized has numerous constraints, Research endeavour goes to develop the algorithm exploited in E-learning |
| Decision trees | A tree like sculpt of assessments along with their consequences | RapidMiner, Smiles,Simple Decision Tree, GATree, KNIME, SilverDicisions, Gambit, YaDT | Assortment of attributes, data training, elucidation of data | A minute amend in dataset would fetch a vast alteration in decision trees, loss of modernization, thorny to shift because of its nature and size etc |
| Text mining | Lashing information from text | Carrot2, GATE, Gensim, OpenNLP, Orange, Stanbol, KNIME, PLOS, Pub Gene | Penetrating of social media, consumer affiliation management, edification, digital humanities | A lot of gratis text is accessible in data anthology, data is amorphous, possibility of syntactic and semantic faults in data, resource progress is intricate |
| Sequential pattern | Regulating objects in a meticulous progression | Mining co, XAffinity (TM), SPMF | Catastrophe prophecy, Recognition of proceeding pills | Immense storage capacity intended for database, Ridge in a store. |
| Prediction | Assessment based on preceding data | EDM, SPSS, WEKA, R | Climate prophecy, student learning performance | Owing to change in future drift, precedent consequences fall short. |

## VIII. CONCLUSION AND FUTURE WORK

Even though the recent techniques still rely on moderately effortless approaches with restricted competence, reassuring consequences have been attained, and the reimbursements of KDD tools have been persuasively established in the wide series of relevance provinces. The amalgamation of critical and realistic desires and afterward strong research interests permits us also anticipate a potential vigorous growth of the domain, depicting KDD tools into the conventional of business applications. Data mining knowledge is a relevance oriented expertise. It is not merely an easy exploration, inquiry or conveys on a meticulous database, but as well evaluates, amalgamates and reasons these facts to direct the elucidation of realistic tribulations and discover the association among proceedings, and yet to envisage potential behaviour through utilizing the accessible data. Data mining has extensive relevance domain roughly in all diligence where the data is produced that's why data mining is deemed as one of the most significant frontiers in database and information organizations and is considered as one of the mainly capable interdisciplinary advancement in Information Technology. Nevertheless seclusion, refuge and exploitation of information are the immense predicaments if it is not concentrated correctly; data mining conveys an assortment of profits to production, society, governments as well as individual. Optimistically diverse classifications and clustering algorithms and its significances would be reviewed in future work.

### REFERENCES

[1] Kaithekuzhical Leena Kurien and Dr. Ajeet Chikkamannur, "A Survey on Methodology of Fraud Detection using Data Mining", International Journal of Trend in Scientific Research and development, Vol.1 ,Issue.6, pp. 38-42, 2017.

[2] D. Ramesh, B. Vishnu Vardhan, "Data Mining Techniques and Applications to Agricultural Yield Data", International Journal of Advanced Research in Computer and Communication Engineering Vol.2, Issue.9, September 2013.

[3] Fayyad. U, Piatetsky-Shapiro. G., and Smyth. P, "From Data Mining to Knowledge Discovery: An Overview", MIT Press, pp.1-36,1996.

[4] H. Benjamin Fredrick David and A. Suruliandi, "Survey on Crime analysis and Prediction using Data mining techniques", ICTACT Journal on Soft Computing, Vol.7, Issue.3, pp.1459-1466, 2017.

[5] R. Tamilsevi and S. Kalaiselvi, "An Overview of Data Mining Techniques and Applications", International Journal of Science and Research, Vol.2, Issue.2, pp. 506-509, 2013.

[6] "Data mining tools", by Ralf Mikut, Markus Reischl, 2011, "Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery".

[7] Er. Rimmy Chuchra, "Use of Data Mining Techniques for the Evaluation of Student Performance: A Case Study" International Journal of Computer Science and Management Research Vol.1, Issue.3, October 2012.

[8] K. Chitra Lekha and Dr. S. Prakasam, "Data mining Techniques in detecting and predicting Cyber crimes in Banking sector", IEEE –

[9] Aakanksha Bhatnagar, Shweta P. Jadye, Madan Mohan Nagar" Data Mining Techniques & Distinct Applications: A Literature Review" International Journal of Engineering Research & Technology (IJERT) Vol.1 Issue 9, November- 2012.

[10] Ruxandra-Ştefania PETRE, "Data mining in Cloud Computing" Database Systems Journal vol. III, no. 3/2012.

[11] K. Chitra Lekha and Dr. S. Prakasam, "Implementation of Data mining techniques for Cyber crime Detection", International Journal of Engineering, Science and Mathematics, Vol.7, Issue.4, pp. 607-613, 2018.

[12] K.Chitra Lekha and Dr. S. Prakasam, "A Survey on Data Mining Techniques in Cyber Crime", IEEE –International Conference on Electrical, Electronics, Computers, Communications, Mechanical and Computing, No.2, January 2018.

[13] Mr. S. P. Deshpande and Dr. V. M. Thakare "Data Mining System and Applications: A Review" International Journal of Distributed and Parallel systems, Vol.1, Issue1, September 2010.

[14] Dr. E. Kesavulu Reddy, "Recent Trends in Data mining Techniques", International Journal of Advance research in Computer Science and Management Studies, Vol.3, Issue.9, pp. 254-262, 2015.

[15] H.-P. Kriegel, K. M. Borgwardt, P. Kröger, A. Pryakhin, M. Schubert, and A. Zimek, "Future trends in data mining," Data Mining Knowledge Discovery,Vol.15, Issue.1, pp.87-97, 2007.

### Authors Profile

K. Chitra Lekha pursuing PhD as full time research scholar in Department of Computer Science and Application, SCSVMV University, Enathur. She has received her M.Phil (Computer Science) from SCSVMV University in 2016, M.Sc (Computer Science) from Annamalai University in 2015 and B.Sc (Computer Science) from Annamalai University in 2013. She has presented 2 papers in IEEE Conference and in a National level Conference and published 4 papers in International Journals. Her research interest includes Data mining, Machine Learning, Software Engineering and Cyber Security.

Dr.S.Prakasam received his PhD from SCSVMVUniversity, Enathur, in 2011, M.C.A degree from Madras University in 2000, B.Sc (Mathematics) from Madras University in 1996. He is presently working as Associate Professor in Department of Computer Science and Application, Sri Chandrasekharendra Saraswathi Viswa Mahavidyalaya (SCSVMV) University, Enathur. He possesses 16 years of vast experience in Computer Science and Applications and has guided many research scholars. He has presented 3 papers in International Conference and he is the author of 40 papers in International journals of repute. His research interests include Knowledge Engineering, Software Agents, and Data Mining.