

## **Data Mining Between Classical and Modern Applications: A Review**

*Ammar Th. Yaseen*

*ammarthaheer@yahoo.com*

*Department of Computer Science,  
College of Computer Science and Mathematics,  
University of Mosul, Mosul, Iraq*

**Received on: 20/04/2021**

**Accepted on: 02/06/2021**

### **ABSTRACT**

Data mining (DM) is an incredible innovation with extraordinary potential to help organizations centre around the main data in the information they have gathered about the conduct of their clients and likely clients. It finds data inside the information that inquiries and reports can't viably uncover. Overall, DM (to a great extent called information or data revelation) is the route toward analysing data according to substitute perspectives and summarizing it into significant information - information that can be used to assemble pay, diminishes costs, or both. DM writing computer programs is one of different logical gadgets for separating data. It grants customers to separate data from a wide scope of estimations or focuses, organize it, and summarize the associations perceived. In reality, DM is the path toward finding associations or models among numerous fields in enormous social datasets. Procedures used in DM measure come from a mix of computational strategies including Artificial Intelligence (AI), Statistics, Machine Learning (ML), and Database (DB) Systems. Aside from the centre techniques used to do the investigation, the cycle of DM can include different pre-handling ventures preceding executing the mining method. Also, a post-preparing stage is normally utilized to picture the aftereffects of the investigation (for example perceived examples or recovered data) in an instinctive and simple to-impart way. From a wide perspective, there are two significant standards of methods: expectation and information disclosure. It includes four sub-groups: a) Classification, Prediction and Regression, b) Clustering, c) Association Rule and Sequence Pattern Mining, and d) Outliers and Anomaly Detection. What's more, there are some generally new and energizing zones of information investigation, for example, spatial DM and graph DM that have been made conceivable through the structure squares of DM techniques. This survey not just advantages analyst to create solid examination subjects and distinguish gaps in the research areas yet additionally helps experts for data mining and Big Data (BD) software framework advancement.

**Keywords:** knowledge discovery in databases; data mining; big data;

### **1. Introduction**

The knowledge discovery in database (KDD) is noted with advancement of strategies and procedures for utilizing information. Quite possibly the main strides of the KDD is the DM. DM is the interaction of example disclosure and extraction where tremendous measure of information is included [1]. Information size is for the most part developing from one day to another. The need to see colossal, complex, information progressed datasets has now extended by and large the changed fields of advancement, business and science. With this tremendous proportion of data, the ability to isolate important data concealed in this gigantic proportion of data and to circle back to the

data is getting logically critical in the present genuine world. The route toward applying PC based information system, including new techniques, for discovering information from data is called DM [2] [3].

DM with BD has been generally utilized in the lifecycle of soft items that rate from the plan with creation steps to the assistance step. A thorough investigation of DM and BD with a survey of its software in the phases of its stages won't just profit specialists to create solid examination subjects and distinguish gaps of research area yet in addition help experts for DM software framework improvement [4].

In this research, a concise explanation of DM associated points is introduced firstly. A diagram of DM and principle substance of the diagram stages are produced in which ordinarily utilized information readiness and pre-handling draws near, DM capacities and procedures, and exhibitions pointers are summed up. At that point, a far reaching audit covering numerous articles on DM or BD applications in the gadgets business is given by the flowchart according to different perspectives, for example, data handling, DM software, or BD at various steps, and the product utilized in the software. On this premise, outlines of information include various information territories and a system in DM and BD software in the electric business are set up.

## **2. What Can Data Mining Do?**

Regardless of the way that DM is at this point in its soonest arranges, associations in a wide extent of organizations - including retail, account, clinical consideration, transportation, and flying - are as of now using DM instruments and methodologies to abuse recorded data. By using plan affirmation headways and quantifiable and mathematical methods to channel through warehoused information, DM helps specialists with seeing immense real factors, associations, designs, models, exceptional cases and irregularities that may some way or another go concealed. For associations, DM is used to discover models and associations in the data to help make better business decisions [5].

Pattern examination - Reveal the contrast between common clients this month and last. ML techniques are getting renowned bit by bit constantly applications like interference area structure, diabetes mining, email spam gathering, etc. The idea behind the DM is particularly straight it is same like an individual become canny from models and experience. In DM; rules are made by taking the direct of given structure (educational file). By then these norms are used to survey the lead/result for the given conditions [6].

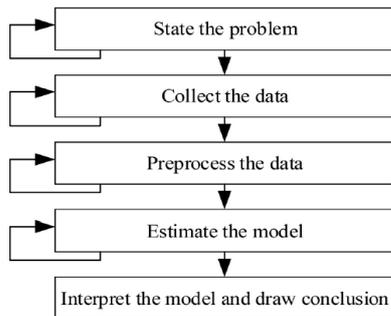
## **3. General Aspects**

DM comprises of extricating data from information put away in DBs to comprehend the information and additionally take choices. Probably the most key DM undertakings are clustering, classification, association rules, outlier detection, and pattern sequences mining. Example mining comprises of finding fascinating, helpful, and unforeseen examples in DBs. This field of exploration has arisen during the 1990s with the original paper of Agrawal and Srikant. That paper presented the Apriori method, intended for finding continuous itemsets, that is gatherings of items (symbols) habitually showing up together in a DB of customer transactions [7].

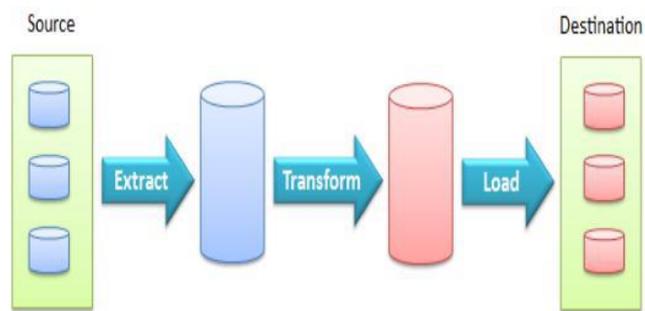
The interest in design mining strategies comes from their capacity to find designs that can be covered up in enormous DBs and that are interpretable by people, and consequently valuable for understanding the information and for decision-making (see Figure 1). For instance, an example of milk; chocolate treats can be utilized to

comprehend client conduct and take vital choices to expand deals, for example, co-advancing items and offering discounts [8].

Information Retrieval (IR) is the gathering of required information resources identified with a data demand from a gathering of information assets. IR investigates words or substance ordering. It is the information on look to information inside text, and furthermore identifying for data of data that express information, and of DBs of pictures or audios, messages. IR frameworks are utilized to lessen data uploaded. For the time of information recovery, data to researched between resulting two ML strategies. The overall methodology information exclude of information recuperation is appeared in Figure 2 [9].



**Figure 1:** The data mining process[10]



**Figure 2:** Extraction data

Successive example mining has some genuine applications since information is encoded as arrangements in numerous fields, for example, bioinformatics, e-learning, market bin investigation, text examination, and page click-stream investigation [7].

Since the Internet of Items (IoT) and progressed data advancements (for instance, radio recurrence recognizable proof labels and smart sensors) are broadly utilized in assembling undertakings for their day by day creation and the board, the item lifecycle the executives measures produce a colossal measure of information. Besides, the amassing of chronicled information in big business asset arranging, inventory network the executives, client relationship the board, and request the executives framework, just as the opportune gathered information by the generally utilized assembling execution framework and disseminated control framework added to the sharp increment of information throughout the long term.

The period of industrial BD has come. “Big Data is high-volume, high-velocity, and high-variety information assets that require new forms of processing to enable enhanced decision making, insight discovery and process optimization”. BD investigation is unequivocally associated with traditional information examination and DM ways to deal with access and interaction these measures of information quick [11]. The electronic health records information actually has a few difficulties that could truly influence or discredit our examination. Instances of such difficulties incorporate circumstances when a patient exits before the investigation closes (alluded to as censoring), patient populaces whose individuals are in generally various conditions of wellbeing and whose results are consequently not practically identical, and an exceptionally factor number of per-patient perceptions [8].

#### **4. Kdd-Knowledge Discovery In Databases**

DM is characterized as a group of conditions, processes, methods that are intended to create significant bits of knowledge, separate examples, and recognize connections from enormous datasets. DM fuses mechanized information extraction, handling, and displaying by methods for a scope of strategies and procedures [12]. DM

projects ordinarily follow an organized cycle or approach as exemplified by Marban and Segovia in 2009, Mariscal and Fernández in 2010. A DM procedure indicates assignments, information sources, yields, and gives rules and directions on how the undertakings are to be executed. Hence, DM procedure gives a group of rules to applying a set of undertakings to accomplish the destinations of a DM research [13].

The establishments of organized DM systems were firstly suggested by Fayyad, Piatetsky-Shapiro and Smyth in 1996, and were at first identified with Knowledge Discovery in Databases (KDD). KDD produces a calculated interaction framework of computed speculations and instruments that help data exclude (knowledge) of information. In KDD, the general way to deal with information disclosure incorporates DM as a particular advance. Accordingly, KDD, with its nine primary advances (showed in Figure 3), has the benefit of thinking about information stockpiling and reach, method covering, translation and perception of outputs, and people PC cooperation. Presentation of KDD additionally formalized more clear differentiation among DM and information examination, with respect to model detailed in Tsai et al. (2015): "...by the data analytics, we mean the whole KDD process, while by the data analysis, we mean the part of data analytics that is aimed at finding the hidden information in the data, such as DM" [14].

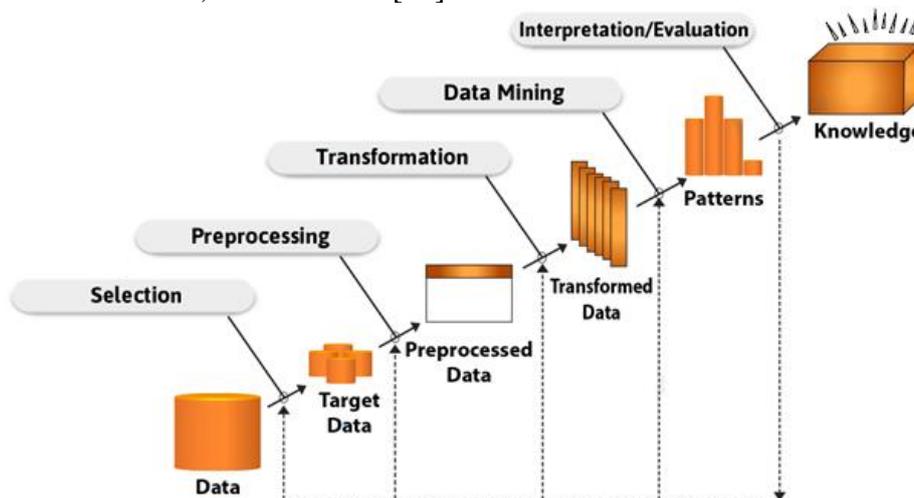


Figure 3: KDD framework

The main nine stages of KDD are as per the following [15]:

**Stage 1:** Learning software domain: It is the initial stage, it is expected to build up a comprehension of the application space and important earlier information attached by distinguishing the objective of KDD cycle of the client's perspective.

**Stage 2:** Dataset generation: Second stage includes choosing a dataset, emphasize on a group of factors or information tests of disclosure is to achieve.

**Stage 3:** Data cleaning and operation: In the 3<sup>rd</sup> stage, fundamental activities to eliminate commotion or exceptions are achieved. Assortment of important data to demonstrate or represent commotion, settling on procedures for dealing with missing information fields, and representing information types, composition, and planning of absent and obscure qualities are additionally thought of.

**Stage 4:** Data reduce and apply: We focusing on identifying helpful properties to address the information, contingent upon the objective of the assignment, use of change techniques to discover ideal properties group for the information is directed.

**Stage 5:** Selecting the capacity of DM: In 5<sup>th</sup> stage, the objective result (e.g., synopsis, clustering, regression, classification) which characterized.

**Stage 6:** Choosing DM method: The 6<sup>th</sup> stage focusing on choosing approach(s) to look for designs in the data, choosing which identify boundaries are proper and coordinating with a specific DM strategy with the general rules of the KDD cycle.

**Stage 7:** DM: In the seventh step, worked by mining the information that is, looking for examples of important in a specific illustrative structure or a group of such display: classification trees, grouping, regression is led.

**Stage 8:** Interpretation: The 8<sup>th</sup> step focus on progression, the excess and unimportant examples are sifted through, pertinent examples are deciphered and imagined in such manner as to make the outcome justifiable to the clients.

**Stage 9:** Using found information: In the last advance, the outcomes are fused with the presentation framework, archived and answered to partners, and utilized as reason for choices.

## **5. CRISP-DM Cross-Industry Standard Process for Data Mining Cycle**

In 2000, as reaction to regular issues and needs, a business-driven approach named Cross-Industry Standard Process for Data Mining (CRISP-DM) acquainted in option with KDD. It likewise merged unique KDD framework and its different expansions. Acknowledgment of DM in numerous life regions prompted CRISP-DM, which is currently the main true norm for DM software. CRISP-DM cycle comprises in six stages [15] [16]:

**Stage 1:** Business understanding: It is focal point of initial stage is to acquire a comprehension of the undertaking destinations and prerequisites from industrial viewpoint related to changing over these in DM issue identifications. Introduction of a fundamental intend to accomplish the goals is likewise remembered for this initial step. The business understanding is generally founded on the given journey definitions and data portrayal.

**Stage 2:** Data understanding: This progression starts with focussing on information arrangement and continuously with expecting for getting comfortable with the information, identify data quality demands, identify knowledge of information, and conceivably distinguish and structure speculations. The information understanding depends on gave data and its documentation.

**Stage 3:** Data preparation: In the 3<sup>rd</sup> stage we covers exercises needed to build the last dataset of the underlying raw data. Data readiness undertakings are achieved more than once. The data readiness comprises of data change, exploratory data examination, and highlight designing. Every one of them can be additionally isolated into more modest sub-steps; e.g., include designing comprises of highlight extraction, highlight choice.

**Stage 4:** Modelling stage: In this progression, different demonstrating methods are chosen and applied followed by adjusting their boundaries. Regularly, a few strategies are utilized for a similar DM issue. In the demonstrating stage, different ML methods can be applied with various boundary alignments. The mix among data and boundary inconstancy can prompt broad rehashing of the model train-test-assessment cycle. On the off chance that the information is enormous scope, the demonstrating stage will have tedious and figure escalated prerequisites.

**Stage 5:** Evaluation of the model(s): The 5<sup>th</sup> stage starts in quality viewpoint and afterward, prior to continuing to conclusive framework arrangement, finds out that the framework(s) accomplishes the industrial targets. Toward the finish of the stage, a choice ought to be arrived on the best way to utilize DM outputs. The assessment stage

can be achieved under different measures for exhaustive checking of the ML frameworks to pick the best framework for the organization stage.

**Stage 6:** Deployment stage: In the last advance, the frameworks are conveyed to empower end clients to utilize the information as reason for choices, or backing in the industrials interaction. Regardless of whether the motivation behind the model is to build information on the data, the information acquired must be helped, introduced, and appropriated such that the end-user could use it. Conditional upon the necessities, the sending step could be just about as straight as producing a state or as intricate as carrying out a iterate DM measure. Sending stage, likewise called creation stage, includes utilization of a prepared ML model to misuse its usefulness, just as the making of an information pipeline into production.

## **6. Crisp-Dm And Kdd: Overlapping And Methodologies**

The utilization of end-to-end DM methodologies, for example, CRISP-DM, KDD cycle, and SEMMA has developed considerably during last decade. In any case, little is known regarding how these approaches are utilized practically [15]. The KDD cycle got predominant in industrial and scholarly areas. Likewise, as course of events based development of DM techniques and cycle models shows, the first KDD DM model filled in as reason for different procedures and interaction models, which tended to different gaps and inadequacies of unique KDD measure. These methodologies broadened the underlying KDD system, yet, augmentation degree has fluctuated going from measure rebuilding to finish modified in focus. An instance, Brachman and Anand [17] and more Gertosio and Dussauchoy [18] (type of contextual investigation) acquainted viable changes with the cycle dependent on iterative nature of interaction just as intelligence. The total KDD measure in their view was improved with strengthening errands and the centre was changed to client's perspective (human-focused methodology), featuring choices that should be made by the client over the span of DM measure. Interestingly, Cabena et al. [19] suggested diverse some stages stressing and enumerating information preparing and revelation tasks. Likewise, in a progression of research Anand and Büchner, Buchner et al. [20] introduced extra DM measure ventures by focusing on transformation of DM cycle to actual configuration. They concentrate on cross-deals (whole life-patterns of real time client), in additional fuse of web information disclosure measure (web-based mining). Moreover, Two Crows DM measure framework is continuously begun structure that has characterized the means in an unexpected way, yet is still near original KDD. At long last, SEMMA dependent on KDD, was created by SAS foundation in 2005. It is characterized as a legitimate association of the useful tools of SAS Enterprise Miner of doing the center assignments of DM. Contrasted with KDD, this is vender explicit interaction framework which restricts its software in various conditions. Likewise, it avoids two stages of unique KDD measure ('Learning Application Domain' also 'Utilizing of Discovered Knowledge') that viewed as fundamental for achievement of DM research. As far as reception, new KDD-based proposition got restricted consideration across the academic and industry world. Accordingly, a large portion of these approaches combined into the CRISP-DM methodology [15].

The following table summarise CRISP-DM and KDD methodologies [13][15]:

**Table 1:** Principle of current data mining operation frameworks and approaches

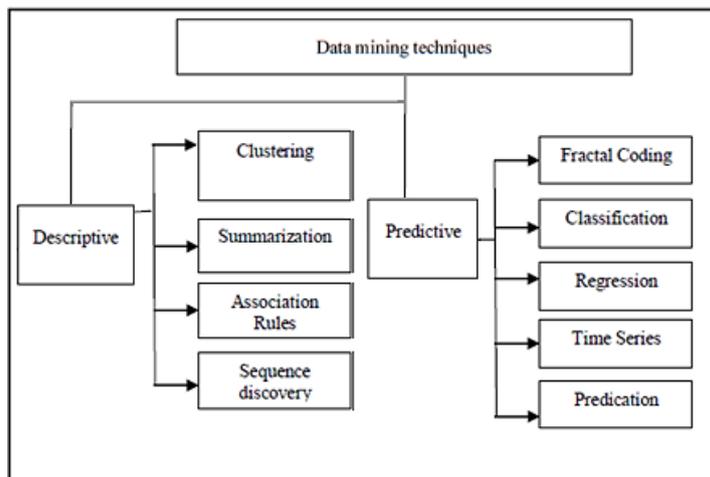
Name	Origin	Basis	Principle	Year
Human-Centred	Academy	KDD	Repeated cycle and intelligence (client's perspective and required choices)	1996, 2004
Cabena et al.	Academy	KDD	Concentrate on information preparing and disclosure undertakings	1997
Anand and Buchner	Academy	KDD	Beneficial advances of coordination of web-mining	1998, 1999
Two Crows	Industry	KDD	Apparatus explicit (SAS foundation), end of certain means	2005
5 A's	Industry	Independent	Valuable advances	2003
6 Sigmas	Industry	Independent	Six Sigma enhancement worldview related to DMAIC execution enhancement framework	2003
CRISP-DM	Academy and industry	KDD	Repeated implementation of stages, huge refinements to errands and yields	2000
Cost et al.	Academy	CRISP-DM	Reconciliation DM of information revelation, criticism components, use of got bits of knowledge upheld by innovations	2005
RAMSYS	Academy	CRISP-DM	Reconciliation and collective research perspectives	2001-2002
DMIE	Academy	CRISP-DM	Reconciliation of variation of modern designing space	2001
Marban	Academy	CRISP-DM	Reconciliation of variation of computer programming space	2007
KDD	Academy and industry	Independent	Device explicit referencing process	2001
ASUM	Industry	CRISP-DM	Device explicit, mix in customary CRISP-DM deft execution method	

## 7. Data Mining Techniques And Frameworks

Lately, PCs and their peripherals have been made less expensive and all the more promptly accessible and in accordance with the improvement of data innovation, different sorts of cutting edge DM methods have touch the business. These modern decade DM methods embrace customary and later complex characterization methods. Both arrangement procedures are for dealing with complex datasets like multidimensionality, client deduction and earlier knowledge, web data, false information focuses that cause over fitting of frameworks, enhancement in people capacity, loud datasets cleaning, mining interactive media datasets and gradual datasets. Overlapped DM methods and methods might be utilized for many the previously announced DBs for predicting the effect and finding significant connections in the information to separate valuable data for information age [2].

Along these lines, assortments of models have been fitted to decide covered up patterns in the data. The methodology that can create the most exact yield and connections design in the noticed datasets is viewed as the most effective in the specific model. Such methodology satisfies the goal of DM. Current DM rehearses uses a scope of model capacities including grouping, regression, classification, finding association rules and arrangement investigation. Be that as it may, the test increments as the interest in DM develops quickly. To deal with these issues without utilizing the conventional statistical strategies, soft processing has arisen to be one of the empowering DM approaches in this field [21].

The cycles in DM are arranged into distinct and prescient (Figure 4). Descriptive mining process give the overall information features in the DB. For predictive mining undertakings, surmising is created of the information for expectations where predict is created on express qualities dependent on designs recognized by ideal outcomes. Descriptive DM, except have any previous goal, gives attributes and portrayals to the dataset [22].



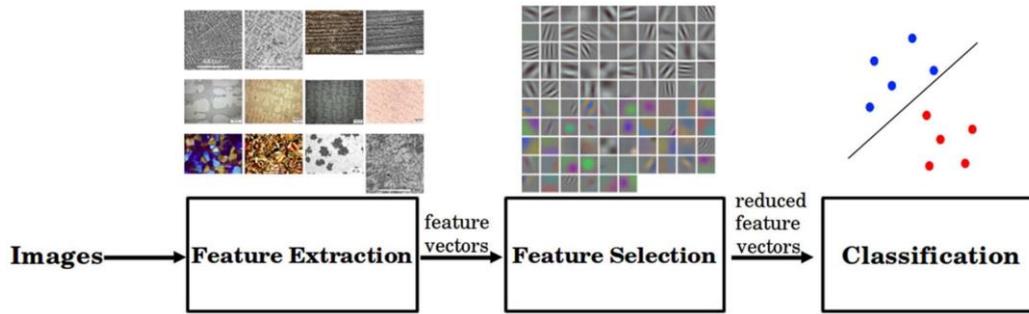
**Figure 4:** Data mining methods

A couple of focus techniques that are used in DM portray the sort of mining and data recovery movement. Tragically, the different associations and courses of action don't for the most part share terms, which can add to the chaos and clear complexity. How about we see some key strategies and instances of how to utilize various devices to build the DM [23].

**A. Association Rules:** Association rules are standard which proposes certain affiliation associations among a bunch of things, (for instance, "happen together" or "one construes the other") in a DB. Given a bunch of exchanges, where each exchange is a bunch of literals (called things), an affiliation rule is a surge of the construction  $X, Y$  ; where  $X$  and  $Y$  are sets of things. The regular meaning of such a standard is that trades of the DB which contain  $X$  will overall contain  $Y$ . A delineation of an affiliation rule is: "30% of ranchers that create wheat also create beats; 2% of all ranchers become both of these things". Here 30% is known as the conviction of the norm, and 2% the assistance of the norm. The issue is to find all affiliation decides that satisfy customer demonstrated least help and least certainty conditions.

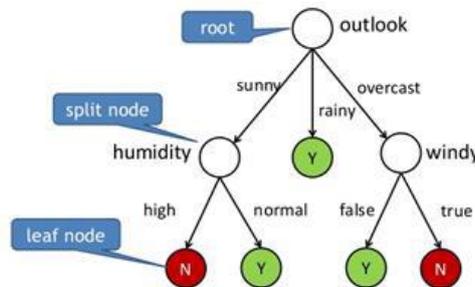
Association (or connection) is well-known and by and large conspicuous and direct DM strategy. Here, you improve on an association between at any rate two things, consistently of comparable sort to recognize plans. For example, when following people's buying affinities, you may separate that a customer reliably buys cream when they buy strawberries, and along these lines recommend that the accompanying time that they buy strawberries they may in like manner need to buy cream [24].

**B. Classification:** You can go through arrangement to assemble a thought of the kind of client, thing, or item by depicting various ascribes to recognize a specific class. For example, you can without a very remarkable stretch gathering vehicles into different sorts (vehicle car, 4x4, convertible) by perceiving different attributes (number of seats, vehicle shape, driven wheels). Given another vehicle, you may apply it into a particular class by differentiating the credits and our known definition [19] (see Figure 5).



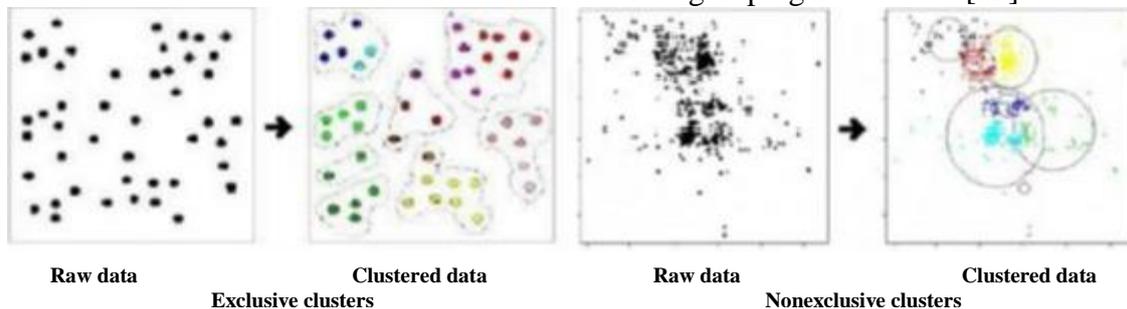
**FIGURE 5:** Schematic explanation of a micrograph classification method on various material composed of feature extraction, feature selection and classification approach

Classification tree is a famous strategy that is utilized to arrange a ward qualified variable dependent on components of at least one forecaster factors. The result is a tree of nodes in connections among the nodes which might be perused to formulate if-then condition (see Figure 6) [23].



**Figure 6:** Classification tree

**C. Clustering (or Grouping):** By taking a gander in any event one attributes or classes, you can assemble solitary pieces of data to outline a development appraisal. At a clear level, gathering is using in any event one credits as your justification perceiving a bunch of relating results. Bunching is useful to recognize particular information since it compares with various models so you can see where the likenesses and scopes agree. Gathering can work the two distinct ways. You can acknowledge that there is a set at explicit point and subsequently use our ID guidelines to check whether you are correct. Grouping is a method gathering the comparable information .In grouping K-NN procedure is finds the distance register and focuses in the chronicled data. Figure 7 shows an illustration of exclusive and non-exclusive grouping of raw data [25].



**Figure 7:** Clustering

**D. Prediction:** Any expectation can be considered as characterization or assessment. The thing that matters is one of accentuation. At the point when DM is utilized to group a telephone line as principally utilized for web access or a MasterCard exchange as deceitful, we don't anticipate having the option to return later to check whether the characterization was right. Our classification might be right or inaccurate, however the

vulnerability is because of fragmented information just: out in reality, the applicable moves have effectively made spot. The telephone is or isn't utilized fundamentally to dial the local ISP. The MasterCard exchange is or isn't fraudulent. With enough endeavours, it is feasible to check [19].

Predictive undertakings feel distinctive in light of the fact that the records are ordered by some anticipated future conduct or assessed future worth. With expectation, the best way to check the exactness of the arrangement is to keep a watch out. Instances of forecast undertakings followed [23]:

Predicting the size in equilibrium which shall be moved if a MasterCard probability acknowledges an equilibrium move offer? Forecasting which clients will left inside next a half year? Forecasting that phone supporters will arrange a worth added administration, for example, three-way calling or audio message.

Methods utilized for arrangement and assessment might be embraced to utilize the expectation by utilizing preparing models when the estimation of the changable to be anticipated is now known, alongside authentic information for those models. The authentic information is utilized to build a model that clarifies the current noticed conduct. At the point where this framework is exercised to present data sources, the outcome is an expectation of further conduct [23].

**E. Statistics:** In this issue of extracting information of information which handled by analysts, well in front of the main AI researches were distributed. For instance, connection examination applies measurable instruments for investigating the relationship between's at least two variables. Grouping investigation offers strategies for finding groups in enormous arrangement of items depicted by vector of qualities. Factor investigation attempts to point the main factors depicting groups. A portion of the well-known strategies that are utilized for managed order undertakings are Linear Discriminants, Quadratic Discriminants, K- nearest Neighbour, Naïve Bays, Logistic Regression and CART [26].

**F. Machine Learning:** Statistical strategies experience issues consolidating abstract, non-quantifiable data in their models. They likewise need to accept different disseminations of boundaries and autonomy of attributes. Different examinations have inferred that ML produces equivalent (and frequently better) predictive accuracy. Its great presentation when contrasted with factual techniques could be ascribed of way which liberated of variable and primary suppositions that describe measurable strategies. Next, shortcoming of mathematical ways to deal with data examination is the issue of deciphering the outcomes [16].

ML and Deep Learning (DL) are parts of the AI. The fundamental objective of ML is to focus on different constant direct insight, in control to display for patterns in data and make improved choices in the possibility dependent on the models. Fundamentally ML order in 4 distinct manners. i. Supervised ,ii. Unsupervised, iii. Semi-supervised iv. Reinforcement.

Administered adapting again separates into 2 classifications of methods (see Figure 8):

Classification: it is the point at which the yield of a changeable is a gathering, for example, "Male" or "Female" and "Human" or "Animal" [27].

Regression: it is the point at which the result variable is genuine worth, for example, "dollars" or "weight".

Unsupervised learning group into 2 classes of methods (see Figure 9): [28]

Clustering: it is only gathering the comparable information. For instance, in the above model every client is placed into one gathering out of the 10 gatherings.

Association: it is "connection between individuals that purchase X additionally will in general purchase Y" [29].

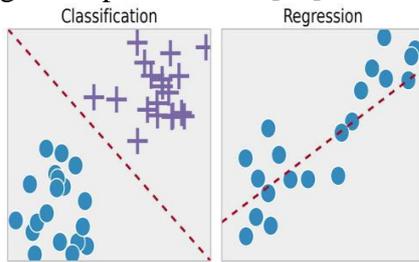


Figure 8: Classification & Regression

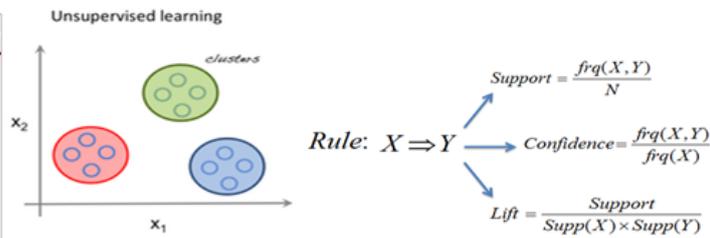


Figure 9: Clustering and association

Semi-supervised learning: Essentially considers the difficulty of when just a little subgroup of the comments has comparing class marks. It utilizes a limited quantity of marked supporting a bigger arrangement of unlabelled data (see Figure 10) [23].

Some of ML approaches are described below.

**G. Neural Networks (NN):** Artificial neural networks include computational frameworks made out of numerous nonlinear processing components ordered in a pattern like natural neuron organizations. A run of the NN has an association esteem related with every node and a weight esteem related with every association. An initiation work administers the terminating of nodes and the spread of data record via network associations in gigantic parallelism. Organization can likewise be prepared with models through association weight changes. NN is a group of techniques and practically identical engineering of different creature, human minds. The framework includes Input, Hidden, Output layers. Each layer determines a specific weight. Data is submitted to the information node, and using the connection of experimentation, the strategy change the loads until it arrive at a specific consummation measures (see Figure 11) [30].

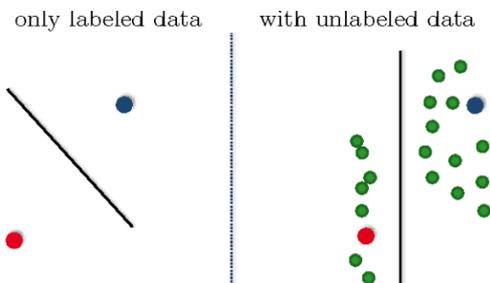


Figure 10: Semi-supervised learning

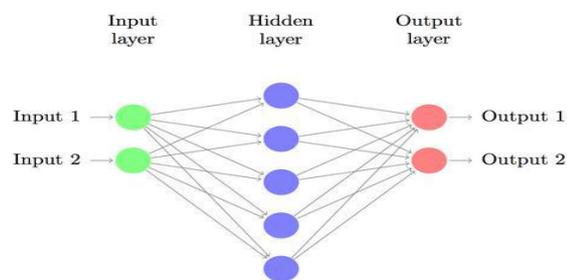


Figure 11: Neural network

**H. Genetic Algorithms (GA)** GAs are search methods dependent on mechanics of natural choice and common genetic activities. They connect natural chose among related structures with an organized at this point randomized record to form a search method with part of imaginative style of human search. Every generation, another group of strings is used pieces and some of the fittest previous generation; an individual new portion is gone after for great measure. While randomized, GA is no randomly basic steps. They productively abuse recorded information to conjecture on added pursuit emphases on anticipated enhancement implementation. A basic GA that yields great outcome, is made out of three administrators to be specific generation, crossover and mutation. It contrasts in every typical enhancement and looking in four methods: GAs use coding boundary set, no simply the boundary. GAs look from a populace of focuses, not a solitary position. GAs utilize objective function data, no subordinates or other

assistant information. GAs utilize likelihood temporary guidelines, not deterministic standards [31].

**I. Support Vector Machines (SVM):** SVMs are the learning approaches which could implement twofold classification and regression assessment assignments. These approaches are getting progressively famous as another worldview of arrangement and learning due to two significant elements. To start with, in contrast to the next classification strategies, SVMs limit the normal mistake as opposed to limiting the classification mistake. Second, SVMs utilize the duality hypothesis of numerical programming to gane a double issue that concedes proficient computing techniques [32].

**J. Fuzzy Logic (FL):** FL, which might be seen as an augmentation of old style sensible frameworks, gives a compelling reasonable structure to managing the issue of information portrayal in a climate of vulnerability and imprecision. A portion of the fundamental qualities of FL identify with the accompanying [33]:

In FL, definite thinking is seen as a restricting instance of rough thinking? In FL everything involves degree? Any intelligent framework can be fuzzed? In FL, information is deciphered as an assortment of flexible or proportionately, fuzzy requirement on an assortment of factors.

**K. Rough Sets (RS) Approaches:** RS theory treats with manages estimate of groups or ideas by methods for paired associations developed from experimental information dependent on the thought of disjointedness and the powerlessness to recognize records. This approximation might assigned to shape frameworks of our objective ideas, and subsequently in its commonplace use, errors under the granular perspective to display development. RS applications to DM by and large continue along the accompanying directions [34]:

Choice standard enlistment from quality worth table? Information filtration by layout age - This primarily includes separating rudimentary squares from information dependent on equality connection. GAs are additionally here and there utilized in this stage for searching.

There are many DM frameworks which include Weka, RapidMiner, Shogun, Scikit-Learn, LibSVM, LibLinear, Vowpal Wabbit, XGBoost, Jupyter notebook, Zeppelin, Kibana, Grafana, Tableau, MATLAB, SAS, R, Python, NumPy, SciPy, PyTorch, Pandas, NLTK, TensorFlow, Keras, Microsoft (CNTK), Caffe, Torch, Apache MXNet, Chainer, Theano, DL4J, Apache Spark MLlib, Spark ML, H2O, FlinkML, Oryx2, KNIME. See [16] for more information.

## 8. Classical Data Mining Applications

Table 2 summarizes the main classical applications in the field of DM [18][20]:

**Table 2: Classical data mining application**

Year	Area	Description
<b>Bayesian applications</b>		
1996	GIS	Surveying the danger of desertification in trouble woods zones
1999	DB System	Diminishing computational time for mining huge DBs
2000	KDD	Dissecting ticket deals information
2004	Medical	Removing human appendage districts
2006	Medical	Building a programmed interracial spikes identification model
2008	Medical	Appraisal of individual danger of supplant or movement in patients determined to have cerebrum tumours going through RT postoperatively
2016	Renewable Energy	Arrangement of mists to decide environmental boundaries prior to building sun based plants
<b>Rule-based applications</b>		
2008	Radiology	Building up a PC helped indicative plan for lung knob location to help radiologists in the identification of cellular breakdown in the lungs from flimsy area processed tomography (CT)

		pictures
2015	Stock trading	Making a few principle pools with a standard based transformative method
2015	Computer security	Utilizing a class association rule-based framework to distinguish network assaults
2016	Medical	Grouping organic information to anticipate the precision of DNA variations and present human interpretable standards
<b>Decision tree applications</b>		
1998	Rule generation	Utilizing choice trees to create fuzzy guidelines for constant esteemed information sources and yields
1999	Learning rules	Joining choice trees with conventional programming to improve learning rules in neural organizations
2002	Identify I/O sets	Applying C4.5 to the outcomes from support learning fused with fuzzy derivation to determine an appropriate activity choice tree
2006	Pattern recognition	Utilizing choice standards for penmanship design order
2008	Rainfall intensity	Building up a heterogeneous various levelled classifier to distinguish downpour regions and precipitation power in those zones
2009	Product impression design	Utilizing survey results to build up a choice tree for client impressions
2009	Image recognition	Building up a gathering choice tree and producing an on-line Random Forest for picture acknowledgment
2012	Artwork	Utilizing C4.5 to characterize the verification of show-stoppers
<b>Neural network applications</b>		
1998	Information retrieval	Creating a word recurrence measure from reports
1999	Speed control	Controlling the speed of a versatile drive framework
2003	Financial	Building up some parametric to direct a neural organization of an objective evaluating recipe
2007	Medical	Changing a neural organization model utilizing acquainted cosmology to make relations between upgrade words and related words
2011	Medical	Utilizing a neural organization for determination of liver disease from clinical pictures
2014	Weather forecasting	Applying diverse back spread to foresee climate conditions
<b>Support vector machines applications</b>		
2002	Industrial	Extricating choice standards for Polycythemia Vera by utilizing a SVM as a classifier
2005	Marketing	Assessing an immediate advertising mailing effort
2007	Medical	Ordering liver infection
2008	Unbalanced data	Building up a weighted consonant mean method to improve the SVM classifier
2009	Human resource	Utilizing a weighted SVM to help dynamic in human asset choice
2012	Recognition	Consolidating a SVM and an internet learning method for hand shape acknowledgment
2014	Financial	Improving the expectation of financial exchange conduct with a one-against-one SVM
<b>K-nearest neighbour applications</b>		
2000	Medical fraud	Arranging remedy composed by General Practitioners as suitable or improper
2011	Network security	Recognizing huge scope assaults like DoS continuously
2016	Big Data	Utilizing K-NN in the planning stage from various preparing information parts from BD
2016	Mechanical	Recognizing distinctive stuff break levels
<b>Case-based reasoning applications</b>		
1997	Legal	Utilizing case rules to decipher court decisions
2002	Medical	Utilizing the significance of features for analysis of lung pathologies
2005	Paper machine	Brushing etymological conditions and fuzzy rationale in the event that based thinking to store arranging and related data to cause arranging a lot quicker when clients to have comparable
2010	Web service	Applying case-based thinking to store arranging and related data to cause arranging a lot quicker when clients to have comparative requirements
2015	Biomedical	Uncovering the ontological application in the field of clinical choice help. Case-based thinking retains and re-establishes experience information for tackling comparable issues with the assistance of the coordinating with approach and characterized interfaces of ontologies
2016	Bankruptcy	Growing new case-based thinking for keeping away from chapter 11 in marginal circumstances
<b>Genetic algorithms applications</b>		
1998	Placement time	Applying a GA to looking for tape feeder course of action to improve arrangement time for chip mounting
2006	Work scheduling	Tracking down the ideal answer for the every minute of every day plan for café
2007	Algorithm	Utilizing a dividing technique to adjust GA to track down a viable arrangement
2009	Medical	Coordinating with lung pictures in numerous subareas
2016	Recognition	Improving execution of human activity acknowledgment with GA and convolutional neural organizations
<b>Rough set applications</b>		
2000	Feature selection	Utilizing harsh sets to give area information communicated to information examination
2004	Applying theory	Applying harsh set hypothesis to choose credits as nodes of a choice tree
2008	Utility program	Characterizing the mathematical qualities in utility projects
2009	Pattern classification	Improving in general execution of unaided ANN

2011	Prediction model	Examining time arrangement information of stunt shrewd value changes
2014	Imperfect data	Dealing with imprecision and vulnerability in information
<b>Fuzzy set applications</b>		
2000	System architecture	Utilizing type-2 fuzzy sets to investigations framework engineering hyper-charts
2001	Transportation	Incorporation fuzzy sets and fuzzy DEA to figure transient creation effectiveness of a city transport organization
2006	Language	Managing word implications in questionable language
2016	Linguistic	Applying fuzzy guidelines to adapt to unbalanced fuzzy sets
<b>Time series applications</b>		
1996	Economic	Inspecting since a long time ago run non-partisanship in 10 distinct nations
1998	Biomedicine	Demonstrating and estimating month to month patient volume
2004	Energy	Guaging hourly power spot costs
2006	Weather	Guaging climate thickness
2010	Rainfall	Guaging the late spring storm precipitation over India
2015	Clinic	Researching singular customer progress and deciding if side effects have improved
<b>Multivariate time series applications</b>		
2000	Environmental	Discovering vegetation and environment variety patterns over the United State
2002	Chemistry	Figuring the water driven boundaries of a stream-spring framework utilizing siphoning test information
2005	Neuroscience	Breaking down multichannel EEGs/MEGs
2007	Medical	Learning and building a conduct profile of an individual from their every day exercises
2015	Agricultural	Guaging transient water system requests
2016	Traffic	Determining traffic framework
<b>Nonparametric applications</b>		
2001	Economics	Building up a few piece based consistency trial of a theory of additivity in nonparametric relapse
2004	Biosystem	Building a quality organization from time arrangement microarray quality articulation information
2011	Statistical Planning	Building up the hypothesis of nonparametric relapse for the traditional instance of reactions missing indiscriminately
2014	Energy	Making on the web condition checking of a breeze power framework
2015	Geoscience	Assessing spatial-transient mountain glacial masses from satellite pictures
2016	Economic	Foreseeing microeconomics and financials toward value changes
<b>Robust regression applications</b>		
1999	Dairy science	Creating transformation conditions for creation, type and wellbeing characteristics
2000	Data analysis	Creating philosophy to maintain a strategic distance from numerous attributions for missing information
2003	Biochemistry	Examining the limiting of the methylphenazinium with twofold abandoned DNA
2006	Chemistry	Identifying anomaly information in logical science
2015	Nero-image	Examining huge neuroimaging partners, recognizing imaging hereditary qualities and staying away from bogus encouraging points in a huge scope examination of mind conduct connections
2016	Data analysis with improving outliers handling	Proposing another powerful relapse to manage case-wise and cell-wise exceptions
<b>Ridge regression applications</b>		
1996	Behaviour	Catching human judgment
2002	Chemometrics	Contrasting adjustment of close to infrared information
2005	Epidemiology	Producing straight models for controlling among metabolites
2007	Hydrology	Applying edge relapse in a component space to figure the hydrologic time arrangement
2016	Computational	Grouping of microarray quality articulation
<b>Nonlinear regression applications</b>		
1998	Aquaculture	Assessing the impacts of dietary phase compound treatment of plant protein slims down for rainbow trout
2003	Economics	Utilizing relapse investigation to break down execution of a college's dial-up modem pool under different time limit approaches and client standards of conduct
2006	Soil	Assessing the effect of peddling, culturing and build up the executives on wheat
2011	Circulatory systems	Diagnosing the disappointment assurance of target frameworks in synthetic plant activity
2014	Environment	Foreseeing the compressive strength of cement
2016	Thermal Biology	Ascertaining and contrasting qualities and limits of metabolic rate and evaporative water misfortune in Australian rodents
<b>Deep learning applications</b>		
2012	Neural Networks	Perceiving German traffic signs
2014	Robotics	Coordinating tactile engine time-arrangement information and the self-association of multimodal intertwined portrayals dependent on a DL approach
2015	Dialogues recognition	Perceiving Chinese discoursed
2016	Human actions	Perceiving human activities in video
<b>K-Means clustering applications</b>		
1998	Image Processing	Division of three dimensional picture information dependent on a novel blend of versatile K-Means setting and information based morphological activities

2005	Physics	Assigning equal progressive N-body collaborations from liquid stream at high Reynolds numbers to gravitational astronomy and sub-atomic elements
2009	Food security	Utilizing seting investigation to feature the linkages between food weakness and neediness
2010	Medical	Changing over the information dark level cerebrum picture to a shading space picture named by grouping
2015	Biomedical	Grouping vectors of perceptions at various time moments and afterward assessing the underlying innervation beat train (IPT)
<b>Fuzzy c-mean applications</b>		
1999	Noise reduction and 3D ultrasonic images	Recreation of pictures from information procured through ultrasonic sensors giving clamor decreased, upgraded pictures
2001	Physics and chemistry of earth	Planning soil information for recognition of contaminated destinations
2008	Medical engineering	Giving a proficient characterization way to deal with computerized isolation of cerebrum MR pictures
2010	Computer in simulations	Improving the location of bogus alerts
2011	Engineering	Giving logical admonition techniques
2015	Mathematical modelling	Lessening warm mistakes of machine instruments in the plan of warm forecast models
2016	Expert systems	Isolating informational indexes by giving a predictable grouping method, which can be utilized for N-dimensional information just as round information
<b>Apriori algorithm applications</b>		
2004	CRM	Finding fuzzy association rules from quantitative exchanges in a store
2006	Database	Removing rules from thick DBs and ward sets of the arrangements of traits in the DB
2007	Database	Applying the RPF-Apriori method to discover fuzzy association rules from fuzzy thing information
2007	Grid technologies	Executing the Apriori method utilizing the Grid administration framework to improve reaction season of the analytics force of different geographic dispersed assets
2012	Web-based	Distinguishing online interruption frameworks to recognize an assortment of assaults and improve the general presentation of discovery frameworks
2016	Medical	Applying association examination to clinical records information to dissect indications and medication mix information
<b>Multidimensional association rules applications</b>		
2001	Weather	Utilizing 2 and m-dimensional settings for mining between value-based association rules to improve prescient capacity
2006	Data warehouse	Characterizing approval rules for clients and protests and wrecking soft data rules to the principle components of a multidimensional model
2009	Web design	Removing helpful data from web clients' entrance ways
2010	Financial	Planning a proficient individual monetary administrations framework by utilizing enlisted data and authentic monetary item data of a client to produce a proposal set
2014	Education	Appling the method to understudies' amalgamation marks assessment framework and breaking down the connected components that impact understudies' blend imprints to assist instructors with improving encouraging strategies and the nature of ability preparing
<b>Quantitative applications</b>		
2005	Dense regions	Catching the qualities of quantitative credits
2011	Biological	Characterizing the interrelationship between qualities
2014	Satellites	Mining all quantitative association rules to improve satellites on-circle execution examination
<b>Sequence discovery applications</b>		
1998	Biology	Finding microbial genome arrangements
2002	Biochemistry	Finding a typical mark of Diacylglyceride kinases, NAD kinases
2008	Biomedical	Removing protein groupings from the protein databank and arranging proteins in the proper overlap classification
2014	Biomedical	Recognizing fleeting connections among drugs and anticipating the following prescriptions to be endorse for a patient

## 9. Modern Data Mining Applications

Table 3 summarizes the main modern applications in the field of DM [2][35] [36]:

**Table 3: Modern Applications**

Seq.	Area	Description
1.	Heart disease	Coronary illness is quite possibly the most widely recognized sicknesses that lead to death in this world. Every year 17.5 million of individuals are passing on because of cardiovascular illness as per World Health Organization reports.
2.	Retrieve hidden information in medical centres	DM is utilized to recover covered up data in clinical focuses that help to anticipate diverse infection. To conquer such issues expectation the event of heart sicknesses utilizing DM procedures and ML methods are assuming fundamental parts for programmed finding of infection in medical care communities.
3.	Large volumes of	Enormous volumes of spatio-transient (ST) information are gathered in a few application spaces,

	spatio-temporal (ST) data	for example, online media, medical services, agribusiness, transportation, and environment science. We momentarily portray the various wellsprings of ST information and the inspiration for examining ST information in various application areas: wrongdoing information, heliophysics, traffic elements, online media, the study of disease transmission/medical care, exactness horticulture, ecological science, neuroscience, environment science.
4.	Medicine science and health informatics	The fields of medication science and wellbeing informatics have gained extraordinary headway as of late and have prompted inside and out investigation that is requested by age, assortment and amassing of huge information. This overview focuses on the survey of on-going investigations utilizing DM and DL approaches for breaking down the particular area information on bioinformatics. It is accepted that in this survey paper, important experiences are accommodated the individuals who are devoted to begin utilizing information investigation techniques in bioinformatics.
5.	Malware detection	DM procedures have been concentrated for malware location in the new Decade. Studies that talking about DM method in dispersed climate.
6.	Sales trends	DM can help spot bargains designs, make more shrewd advancing endeavours, and correctly expect customer dedication. Unequivocal occupations of DM include: Market division - Identify the fundamental ascribes of customers who buy comparable things from your association. Customer foment - Predict which customers are likely going to leave your association and go to a competitor. Deception disclosure - Identify which trades are bound to be bogus. Direct exhibiting - Identify which prospects should be associated with a mailing overview to get the most vital response rate. Intelligent promoting - Predict what every individual getting to a site is no doubt keen on seeing. Market bushel examination - Understand what items or administrations are generally bought together; e.g., brew and diapers.
7.	Bioinformatics	Bioinformatics is a promising zone in the field of solution, biotechnology, drugs plan, microbiology, agribusiness and PC.
8.	Healthcare	DM has been applied effectively in medical services extortion and recognizing misuse cases. Clinical decisions are consistently made reliant on experts' nature and experience instead of on the data rich data concealed in the DB. This preparation prompts bothersome inclinations, goofs and absurd clinical costs which impacts the idea of organization provided for patients. This thought is promising as data exhibiting and assessment gadgets, e.g., DM, can establish a data rich environment which can serve to through and through improve the idea of clinical decisions.. Effective DM applications have given the impulse to the important gatherings to completely used them as they have understood that DM is critical in the procurement of significant data for all areas associated with medical care related businesses.
9.	Medical data	To utilize DM methods to clinical information, specialists' cognizance on the kind of DM methods and their capacities ought to be clear.
10.	Image and video processing	Satellites pictures (for instance, fires, dry spells, crops illnesses, metropolitan turn of events), space (telescope pictures), science picture acknowledgment (for example plant, cells, microscopic organisms), clinical picture acknowledgment (for example attractive reverberation imaging, registered tomography, solography, histopathology), PC vision, programmed picture or sound comments.
11.	Speech	Discourse and language text handling and acknowledgment, discourse acknowledgment, machine interpretation, characteristic language preparing;
12.	Security	Security biometrics verification (for example individuals, faces, step) and network protection (for example organization and asset observing, abnormality discovery, interruption recognition).
13.	Business	Business insight protection, monetary business sectors, stock and swapping scale (for example time-arrangement observing and forecasts).
14.	Robotics	Advanced mechanics and computer games independent route (for example vehicle, drone, plane, submarine), computer games (for example Atari, Dota, Starcraft).
15.	Distributed platforms	At the main sight, general crowd relates BD handling with conveyed stages like Apache Hadoop also Spark. The focusing for Volume in similarity of Veracity trademark just as high-speeds in handling and derivation in foundation.
16.	GPU	The fundamental element of many-centre gas pedals, for example, GPU is their enormously equal engineering permitting them to accelerate methods that include lattice based activities, which are on a basic level of numerous ML/DL executions.
17.	Dataflow	Late progressed DM subjects; i.e., the current interest of handling enormous scope information. The vast majority of the advanced DM instruments depend dataflow models (pipeline or work process). In addition, they coordinated graphical UI (GUI), lean toward an API method.
18.	General applications	Promoting, Customer relationship the executives, Engineering, Medicine investigation, Expert forecast, Web mining, Mobile processing.

### 10. Data Mining And Big Data Challenges

BD is put away a tremendous measure of data, the investigating a lot of data, testing helpful data from different datasets. It is additionally unstructured information, huge size and it is difficult to deal with. There are a few contrasts among DM and BD which summed up in the following table [9]:

**Table 4:** Differences between DM and BD

Area	Data Mining	Big Data
Concept	This refers to strategies to find the important data and fascinating examples with regards to datasets; the datasets could be little records or huge volumes of information	This is a term referring to enormous scope stockpiling and preparing of huge datasets; the dataset becomes quicker than a straightforward DB and past information taking care of models
Value	It is devices and strategies that utilize significant outcomes for choice maker	It is a resource
Processing	Handling alludes to the activity that includes moderately modern inquiry tasks	Processing differs relying upon the capacities of the association overseeing it and applications to measure and examinations the information
Task	Not all DM assignments manage BD	All BD undertakings include DM

The following table summarises the main issues and their challenges in DM and BD [11].

**Table 5:** Main issues and challenges

S.no	The main issue	Major challenges
1.	Helpless information quality. For instance, boisterous information, filthy data	Big DM Platform
2.	Excess information is transferred from different sources	Information Sharing and Data Privacy
3.	Security, protection of the companies	Domain and Application Knowledge
4.	Greater expense, less flexibility	Big DM Algorithm

As individuals compose words with faults, to allow them to compose or look with appropriate syntax and organized words, text mining method is utilized. Text mining implies excluding of the information which isn't comfortable to anybody. In the event that we contrast web searching and text mining, those are incomprehensibly unique in relation to one another. Text mining makes it simple to get significant and organized information from the unpredictable information designs. It is truly not a simple task for the PCs to comprehend the unstructured information and make it organized. Individuals can play out this assignment with no further endeavours because of the accessibility of various linguistic methods. Information accessible in the content arrangement has significantly more significance and that is the reason text mining is producing a lot of business esteem. An examination by expressed that DM addresses the inference of a significant example or standards spatial DB of deciding a specific problems. DM is not quite the same as text mining. An examination by called attention to text mining significantly most perplexing than DM since it includes unpredictable or unstructured information designs, though DM is managing the organized arrangements of information [37].

## 11. Conclusion

This paper has introduced various methods of DM ,BD, and ML approaches. The general objective is to assess the utilization of the DM methods for various sorts of frameworks. Every method has diverse objective and objective to arrange the dataset in various ways. Like clustering based methods has capacity to build up the group of homogeneous information component from a given dataset like sex attribute partition dataset into males and females. Additionally Apriori method for assessing the connections among attributes like reward property has incredible impact on the month to month pay of an individual.

With the quick advancement of different situating procedures like GPS, cell phones and distant detecting, spatio-temporal information has become progressively accessible these days. Mining significant information from spatio-transient information is fundamentally essential to some true applications including human mobility understanding, smart transportation, metropolitan arranging, public security, medical services and natural administration. Associations may go through for investigation of BD to having better choices, in this way BD examination is being focused lately. For

tracking down the covered qualities from BD, society requires new plans or systems. Customary DM methods accept that the information is midway gathered, memory-resident, and static. It is trying to deal with the large-scale information and interaction them with extremely restricted assets. For instance, a lot of information are immediately delivered and put away at different areas. It turns out to be progressively costly to bring together them in a solitary place.

This research might be summed up in the following:

1. Most DL systems are created in world's biggest programming organizations like Google, Facebook, and Microsoft. These organizations have colossal measures of information, elite foundations, human insight and venture assets.
2. BD environments, for example, Apache Flink, Apache Spark, Cloudera Oryx 2 include work in ML libraries of huge scope DM essentially for plain information. These ML libraries are presently in an advancing state however the force of the entire biological system is critical.
3. The pattern shows a high number of intuitive information investigation (like Python) and information representation apparatuses supporting leaders.

We infer that there is still likelihood to improve different methods and procedures for DM.

**REFERENCE**

- [1] C. Bellinger, M. Shazan, M. Jabbar, O. Zaiane, and A. Osornio-vargas, “Open Access A systematic review of data mining and machine learning for air pollution epidemiology,” pp. 1–19, 2017.
- [2] S. A. Lashari, R. Ibrahim, N. Senan, and N. S. A. M. Taujuddin, “Application of Data Mining Techniques for Medical Data Classification: A Review,” vol. 06003, pp. 1–6, 2018.
- [3] A. Al Abd Alazeez, S. Jassim, and H. Du, “EINCKM: An Enhanced Prototype-based Method for Clustering Evolving Data Streams in Big Data,” *Proc. 6th Int. Conf. Pattern Recognit. Appl. Methods*, no. Icpram, pp. 173–183, 2017.
- [4] N. Jothi, N. Aini, A. Rashid, and W. Husain, “Data Mining in Healthcare – A Review,” *Procedia - Procedia Comput. Sci.*, vol. 72, pp. 306–313, 2015.
- [5] A. T. Y. T. A. A. Alazeez, “HPPD: A Hybrid Parallel Framework of Partition-based and Density-based Clustering Algorithms in Data Streams Ammar,” *Raf. J. Comp. Math's.*, vol. 1, no. 1, pp. 67–82, 2020.
- [6] J. L. Lobo, J. Del, S. Eneko, O. Albert, and B. Francisco, “CURIE: A Cellular Automaton for Concept Drift Detection,” *arXiv Prepr. arXiv ...*, vol. 5, no. 3, p. 15, 2020.
- [7] P. Fournier-viger and J. C. Lin, “A Survey of Sequential Pattern Mining,” vol. 1, no. 1, pp. 54–77, 2017.
- [8] P. Yadav, M. Steinbach, V. Kumar, and G. Simon, “Mining Electronic Health Records ( EHRs ): A Survey,” vol. 50, no. 6, pp. 1–40, 2018.
- [9] J. N. Rao and M. Ramesh, “A Review on Data Mining & Big Data , Machine Learning Techniques,” no. 6, pp. 914–916, 2019.
- [10] R. Pruengkarn, K. W. Wong, and C. C. Fung, “A Review of Data Mining Techniques and Applications,” 2016.
- [11] S. Lv and H. Kim, “applied sciences A Review of Data Mining with Big Data towards Its Applications in the Electronics Industry,” no. Dm, pp. 1–34, 2018.
- [12] C. Series, “A review on bioinformatics using data mining techniques A review on bioinformatics using data mining techniques,” 2019.
- [13] R. Venkatesh, K. Chaitanya, T. Bikku, and R. Paturi, “REVIEW ARTICLE A Review on Biomedical Mining,” vol. 15, pp. 629–637, 2019.
- [14] A. Souri and R. Hosseini, “A state - of - the - art survey of malware detection approaches using data mining techniques,” *Human-centric Comput. Inf. Sci.*, 2018.
- [15] V. Plotnikova, M. Dumas, and F. Milani, “Adaptations of data mining methodologies : a systematic literature review,” 2020.
- [16] G. Nguyen *et al.*, “Machine Learning and Deep Learning frameworks and libraries for large-scale data mining: a survey,” *Artif. Intell. Rev.*, vol. 52, no. 1, pp. 77–124, 2019.
- [17] F. A. Rafrastara, G. S. Member, and Q. Deyu, “A Survey and Taxonomy of

- Distributed Data Mining Research Studies : A Systematic Literature Review,” vol. 14, no. 5, 2016.
- [18] P. Sunhare, R. R. Chowdhary, and M. K. Chattopadhyay, “Internet of things and data mining : An application oriented survey,” *J. King Saud Univ. - Comput. Inf. Sci.*, no. xxxx, 2020.
- [19] W. G. Touw *et al.*, “Data mining in the life science swith random forest: A walk in the park or lost in the jungle?,” *Brief. Bioinform.*, vol. 14, no. 3, pp. 315–326, 2012.
- [20] G. Atluri, A. Karpatne, and V. Kumar, “Spatio-Temporal Data Mining : A Survey of Problems,” vol. 51, no. 4, 2018.
- [21] M. Chen, S. Mao, and Y. Liu, “Big Data: A Survey,” *Mob. Networks Appl.*, vol. 19, no. 2, pp. 171–209, Jan. 2014.
- [22] K. Stensbo-Smidt, C. Igel, A. Zirm, and K. S. Pedersen, “Nearest neighbour regression outperforms model-based prediction of specific star formation rate,” *Proc. - 2013 IEEE Int. Conf. Big Data, Big Data 2013*, pp. 141–144, 2013.
- [23] F. E. Bock, R. C. Aydin, C. J. Cyron, N. Huber, S. R. Kalidindi, and B. Klusemann, “A Review of the Application of Machine Learning and Data Mining Approaches in Continuum Materials Mechanics,” vol. 6, no. May, 2019.
- [24] A. T. Y. A. A. Al Abd Alazeez, “AEPRD: An Enhanced Algorithm for Predicting Results of Orthodontic Operations,” *J. Educ. Sci.*, vol. 30, no. 1, pp. 173–190, 2021.
- [25] A. T. Y. Al Abd Alazeez, S. Jassim, and H. Du, “SLDPC : Towards Second Order Learning for Detecting Persistent Clusters in Data Streams,” *2018 10th Comput. Sci. Electron. Eng.*, vol. 978-1-5386, pp. 248–253, 2018.
- [26] A. Patri and Y. Patnaik, “Random Forest and Stochastic Gradient Tree Boosting Based Approach for the Prediction of Airfoil Self-noise,” *Procedia Comput. Sci.*, vol. 46, no. Iccit 2014, pp. 109–121, 2015.
- [27] M. J. Zaki and L. Wong, “Data mining techniques,” *WSPC/Lecture Notes Ser. 9in x 6in*, 2003.
- [28] S. Lertampaiporn, C. Thammarongtham, C. Nukoolkit, B. Kaewkamnerdpong, and M. Ruengjitchatchawalya, “Identification of non-coding RNAs with a new composite feature in the Hybrid Random Forest Ensemble algorithm,” *Nucleic Acids Res.*, vol. 42, no. 11, pp. 1–13, 2014.
- [29] A. Acharjee *et al.*, “Data integration and network reconstruction with ???omics data using Random Forest regression in potato,” *Anal. Chim. Acta*, vol. 705, no. 1–2, pp. 56–63, 2011.
- [30] R. Ben Ali, R. Ejbali, and M. Zaied, “Detection and Classification of Dental Caries in X-ray Images Using Deep Neural Networks,” no. c, pp. 223–227, 2016.
- [31] A. Al Abd Alazeez, S. Jassim, and H. Du, “TPICDS: A Two-phase Parallel Approach for Incremental Clustering of Data Streams,” *24th Int. Eur. Conf. Parallel Distrib. Comput.*, 2018.
- [32] M. Chakarverti, N. Sharma, and R. R. Divivedi, “Prediction Analysis Techniques

- of Data Mining :,” 2019.
- [33] P. Bonissone, J. M. Cadenas, M. Carmen Garrido, and R. Andrés Díaz-Valladares, “A fuzzy random forest,” *Int. J. Approx. Reason.*, vol. 51, no. 7, pp. 729–747, 2010.
- [34] S. Yun, “Research of Big Data Analysis on Rough Set and Electromagnetism-like Mechanism Algorithm,” *IEEE Int. Conf. Comput. Inf. Technol. Res.*, no. 978, pp. 923–926, 2014.
- [35] C. Beyene, “Survey on Prediction and Analysis the Occurrence of Heart Disease Using Data Mining Techniques,” vol. 118, no. 8, pp. 165–174, 2018.
- [36] A. Al Abd Alazeez, S. Jassim, and H. Du, “EDDS: An Enhanced Density-Based Method for Clustering Data Streams,” *2017 46th Int. Conf. Parallel Process. Work.*, pp. 103–112, 2017.
- [37] S. A. Salloum, M. Al-Emran, A. A. Monem, and K. Shaalan, “A survey of text mining in social media: Facebook and Twitter perspectives,” *Adv. Sci. Technol. Eng. Syst.*, vol. 2, no. 1, pp. 127–133, 2017.