

Issues and Techniques in Data Mining and Cluster Techniques

Pankaj Kumar, Paritosh Kumar and Gajraj Singh

*Ramjas College, University of Delhi
pkumar240183@gmail.com
CEM Kapurthala, Punjab(India)
paritosh200623@gmail.com
Ramjas College, University of Delhi
gajraj76@gmail.com*

Abstract

The objective of the paper is to show the issues to be faced in Data mining and various clustering Techniques. Clustering is being widely used in many application including medical, finance and etc. Clustering may be applied on database using various approaches, based upon distance, density, hierarchy, and partition. Data mining is becoming an increasingly important tool to transform this data into information. It is commonly used in a wider range of profiling practices, such as marketing, surveillance, fraud detection and scientific discovery.

Keywords Clustering, Knowledge Discovery, Noise, Data mining,

Introduction

Data mining is the process of extracting hidden patterns from data. As more data is gathered, with the amount of data doubling every three years, [1] data mining is becoming an increasingly important tool to transform this data into information. It is commonly used in a wider range of profiling practices, such as marketing, surveillance, fraud detection and scientific discovery. Data mining, the extraction of hidden predictive information from large databases, is a powerful new technology with great potential to help companies focus on the most important information in their data warehouses. Data mining tools predict future trends and behaviours, allowing businesses to make proactive, knowledge-driven decisions. The automated, prospective analyses offered by data mining move beyond the analyses of past events provided by retrospective tool typical of decision support systems.

Data mining tools can answer business questions that traditionally were too time-consuming to resolve. They scour databases for hidden patterns, finding predictive information that experts may miss because it lies outside their expectations. Most companies already collect and refine massive quantities of data. Data mining techniques can be implemented rapidly on existing software and hardware platforms to enhance the value of existing information resources, and can be integrated with new products and systems as they are brought on-line. When implemented on high performance client/server or parallel processing computers, data mining tools can analyze massive databases to deliver answers to questions such as, "Which clients are most likely to respond to my next promotional mailing, and why?" Data mining is the process of discovering meaningful new correlations, patterns, and trends by sifting through large amounts of data stored in repositories, using pattern recognition technologies as well as statistical and mathematical techniques". However, really data mining turns databases into knowledge bases which is one of the fundamental components of expert systems. Data mining is the use of automated data analysis techniques to uncover previously undetected relationships among data items. Data mining often involves the analysis of data stored in a data warehouse. Three of the major data mining techniques are regression, classification and clustering. Data Mining, also popularly known as Knowledge Discovery in Databases (KDD) [2], refers to the nontrivial extraction of implicit, previously unknown and potentially useful information from data in databases. While data mining and knowledge discovery in databases (or KDD) are frequently treated as synonyms, data mining is actually part of the knowledge discovery process. The following figure (Figure 1. 1) shows data mining as a step in an iterative knowledge discovery process. The Knowledge Discovery in Databases process comprises of a few steps leading from raw data collections to some form of new knowledge. The iterative process consists of the following steps:

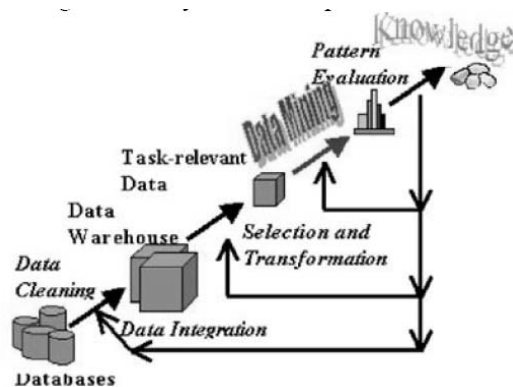


Fig 1: An iterative Knowledge Discovery Process

- a. Data cleaning: also known as data cleansing, it is a phase in which noised data and irrelevant data are removed from the collection.
- b. Data integration: at this stage, multiple data sources, often heterogeneous, may be combined in a common source.

- c. Data selection: at this step, the data relevant to the analysis is decided on and retrieved from the data collection.
- d. Data transformation: also known as data consolidation, it is a phase in which the selected data is transformed into forms appropriate for the mining procedure.
- e. Data mining: it is the crucial step in which clever techniques are applied to extract patterns potentially useful.
- f. Pattern evaluation: in this step, strictly interesting patterns representing knowledge are identified based on given measures.
- g. Knowledge representation: is the final phase in which the discovered knowledge is visually represented to the user. This essential step uses visualization techniques to help users understand and interpret the data mining results.

III Issues in Data Mining

Data mining algorithms embody techniques that have sometimes existed for many years, but have only lately been applied as reliable and scalable tools that time and again outperform older classical statistical methods. While data mining is still in its infancy, it is becoming a trend and ubiquitous. Before data mining develops into a conventional, mature and trusted discipline, many still pending issues have to be addressed. Some of these issues are addressed below. Note that these issues are not exclusive and are not ordered in any way.

Security and social issues: Security is an important issue with any data collection that is shared and/or is intended to be used for strategic decision-making. In addition, when data is collected for customer profiling, user behaviour understanding, correlating personal data with other information, etc., large amounts of sensitive and private information about individuals or companies is gathered and stored. This becomes controversial given the confidential nature of some of this data and the potential illegal access to the information. Moreover, data mining could disclose new implicit knowledge about individuals or groups that could be against privacy policies, especially if there is potential dissemination of discovered information. Another issue that arises from this concern is the appropriate use of data mining. Due to the value of data, databases of all sorts of content are regularly sold, and because of the competitive advantage that can be attained from implicit knowledge discovered, some important information could be withheld, while other information could be widely distributed and used without control.

User interface issues: The knowledge discovered by data mining tools is useful as long as it is interesting, and above all understandable by the user. Good data visualization eases the interpretation of data mining results, as well as helps users better understand their needs. Many data exploratory analysis tasks are significantly facilitated by the ability to see data in an appropriate visual presentation. There are many visualization ideas and proposals for effective data graphical presentation. However, there is still much research to accomplish in order to obtain good visualization tools for large datasets that could be used

to display and manipulate mined knowledge. The major issues related to user interfaces and visualization are "screen real-estate", information rendering, and interaction. Interactivity with the data and data mining results is crucial since it provides means for the user to focus and refine the mining tasks, as well as to picture the discovered knowledge from different angles and at different conceptual levels.

Mining methodology issues: These issues pertain to the data mining approaches applied and their limitations. Topics such as versatility of the mining approaches, the diversity of data available, the dimensionality of the domain, the broad analysis needs (when known), the assessment of the knowledge discovered, the exploitation of background knowledge and metadata, the control and handling of noise in data, etc. are all examples that can dictate mining methodology choices. For instance, it is often desirable

to have different data mining methods available since different approaches may perform differently depending upon the data at hand. Moreover, different approaches may suit and solve user's needs differently. Most algorithms assume the data to be noise-free. This is of course a strong assumption. Most datasets contain exceptions, invalid or incomplete information, etc. which may complicate, if not obscure, the analysis process and in many cases compromise the accuracy of the results. As a consequence, data preprocessing (data cleaning and transformation) becomes vital. It is often seen as lost time, but data cleaning, as time-consuming and frustrating as it may be, is one of the most important phases in the knowledge discovery process. Data mining techniques should be able to handle noise in data or incomplete information. More than the size of data, the size of the search space is even more decisive for data mining techniques. The size of the search space

is often dependent upon the number of dimensions in the domain space. The search space usually grows exponentially when the number of dimensions increases. This is known as the curse of dimensionality. This "curse" affects so badly the performance of some data mining approaches that it is becoming one of the most urgent issues to solve.

Performance issues: Many artificial intelligence and statistical methods exist for data analysis and interpretation. However, these methods were often not designed for the very large data sets data mining is dealing with today. Terabyte sizes are common. This raises the issues of scalability and efficiency of the data mining methods when processing considerably large data. Algorithms with exponential and even medium order polynomial complexity cannot be of practical use for data mining. Linear algorithms are usually the norm.

In some cases, sampling can be used for mining instead of the whole dataset. However, concerns such as completeness and choice of samples may arise. Other topics in the issue of performance are incremental updating, and parallel programming [3]. There is no doubt that parallelism can help solve the size problem if the dataset can be subdivided and the results can be merged later. Incremental updating is important for merging results from parallel mining, or updating data mining results when new data becomes available without having to reanalyze the completed dataset.

Data source issues: There are many issues related to the data sources,

some are practical such as the diversity of data types, while others are philosophical like the data glut problem. We certainly have an excess of data since we already have more data than we can handle and we are still collecting data at an even higher rate. If the spread of database management systems has helped increase the gathering of information, the advent of data mining is certainly encouraging more data harvesting. The current practice is to collect as much data as possible now and process it, or try to process it, later. The concern is whether we are collecting the right data at the appropriate amount, whether we know what we want to do with it, and whether we distinguish between what data is important and what data is insignificant. Regarding the practical issues related to data sources, there is the subject of heterogeneous databases and the focus on diverse complex data types. We are storing different types of data in a variety of repositories. It is difficult to expect a data mining system to effectively and efficiently achieve good mining results on all kinds of data and sources. Different kinds of data and sources may require distinct algorithms and methodologies. Currently, there is a focus on relational databases and data warehouses, but other approaches need to be pioneered for other specific complex data types. A versatile data mining tool, for all sorts of data, may not be realistic.

Moreover, the proliferation of heterogeneous data sources, at structural and semantic levels, poses important challenges not only to the database community but also to the data mining community.

III Data Mining Techniques

The three data mining techniques are:

- a. Decision trees
- b. Neural networks
- c. Clustering

Decision tree

A decision tree [4] (or tree diagram) is a decision support tool that uses a tree-like graph or model of decisions and their possible consequences, including chance event outcomes, resource costs, and utility. Decision trees are commonly used in operations research, specifically in decision analysis, to help identify a strategy most likely to reach a goal. Another use of decision trees is as a descriptive means for calculating conditional probabilities. In data mining and machine learning, a decision tree is a predictive model; that is, a mapping from observations about an item to conclusions about its target value. More descriptive names for such tree models are classification tree (discrete outcome) or regression tree (continuous outcome). In these tree structures, leaves represent classifications and branches represent conjunctions of features that lead to those classifications. The machine learning technique for inducing a decision tree from data is called decision tree learning, or (colloquially) decision trees. In decision analysis, a "decision tree"—and the closely-related influence diagram—

is used as a visual and analytical decision support tool, where the expected values (or expected utility) of competing alternatives are recalculated. A decision tree consists of 3 types of nodes: -

1. Decision nodes - commonly represented by squares.
2. Chance nodes - represented by circles.
3. End nodes - represented by triangles.

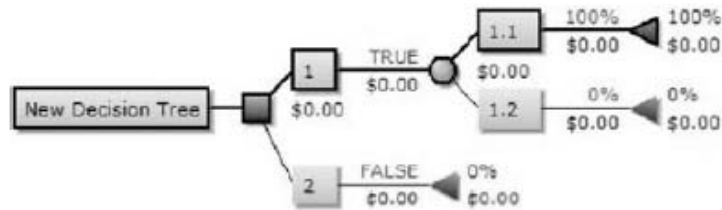


figure 1. 2

Drawn from left to right, a decision tree has only burst nodes (splitting paths) but no sink nodes (converging paths).

Influencediagram

A decision tree can be represented more compactly as an influencediagram, focusing attention on the issues and relationships between events.

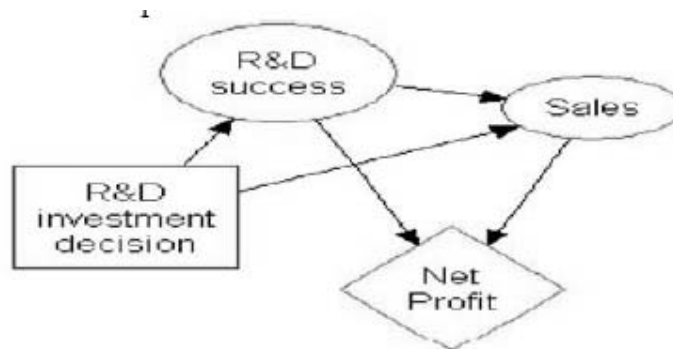


figure 1. 3 influencediagram

Creation of decision nodes

Three popular rules are applied in the automatic creation of classification trees. The Gini rule splits off a single group of as large a size as possible, whereas the entropy and towing rules find multiple groups comprising as close to half the samples as possible. Both algorithms proceed recursively down the tree until stopping criteria are met.

Amongst decision support tools, decision trees (and influencediagrams) have several advantages:

- Decision trees are simple to understand and interpret.
- People are able to understand decision tree models after a brief explanation.
- Have value even with little hard data. Important insights can be generated based on experts describing a situation (its alternatives, probabilities, and costs) and their preferences for outcomes.
- Use a white box model. If a given result is provided by a model, the explanation for the result is easily replicated by simple math. Can be combined with other decision techniques. The following example uses Net Present Value calculations, PERT 3-point estimations (decision #1) and a linear distribution of expected outcomes (decision #2):

b. Neural network

Neural Networks [5] are analytic techniques modeled after the (hypothesized) processes of learning in the cognitive system and the neurological functions of the brain and capable of predicting new observations (on specific variables) from other observations (on the same or other variables) after executing a process of so-called learning from existing data. Neural Networks is one of the Data Mining techniques. The first step is to design a specific network architecture (that includes a specific number of "layers" each consisting of a certain number of "neurons"). The size and structure of the network need to match the nature (e. g., the formal complexity) of the investigated phenomenon. Because the latter is obviously not known very well at this early stage, this task is not easy and often involves multiple "trials and errors." (Now, there is, however, neural network software that applies artificial intelligence techniques to aid in that tedious task and finds "the best" network architecture). The new network is then subjected to the process of "training." In that phase, neurons apply an iterative process to the number of inputs (variables) to adjust the weights of the network in order to optimally predict (in traditional terms one could say, find a "fit" to) the sample data on which the "training" is performed. After the phase of learning from an existing dataset, the new network is ready and it can then be used to generate predictions. Neural network had been used to refer to a network or circuit of biological neurons. The modern usage of the term often refers to artificial neural networks, which are composed of artificial neurons or nodes. Thus the term has two distinct usages: Biological neural networks are made up of real biological neurons that are connected or functionally related in the peripheral nervous system or the central nervous system. In the field of neuroscience, they are often identified as groups of neurons that perform a specific physiological function in laboratory analysis. Artificial neural networks are made up of interconnecting artificial neurons (programming constructs that mimic the properties of biological neurons). Artificial neural networks may either be used to gain an understanding of biological neural networks, or for solving artificial intelligence problems without necessarily creating a model of a real biological system. Thereal, biological nervous system is highly complex and includes some features that may seem superfluous based on an understanding of

artificial networks. In general a biological neural network is composed of a group or groups of chemically connected or functionally associated neurons. A single neuron may be connected to many other neurons and the total number of neurons and connections in a network may be extensive. Connections, called synapses, are usually formed from axon to dendrites, though dendrodendritic microcircuits [1] and other connections are possible. Apart from the electrical signaling, there are other forms of signaling that arise from neurotransmitter diffusion, which have an effect on electrical signaling. As such, neural networks are extremely complex. Artificial intelligence and cognitive modeling try to simulate some properties of neural networks. While similar in their techniques, the former has the aim of solving particular tasks, while the latter aims to build mathematical models of biological neural systems. In the artificial intelligence field, artificial neural networks have been applied successfully to speech recognition, image analysis and adaptive control, in order to construct software agents (in computer and video games) or autonomous robots. Most of the currently employed artificial neural networks for artificial intelligence are based on statistical estimation, optimization and control theory. The cognitive modeling field involves the physical or mathematical modeling of the behaviour of neural systems; ranging from the individual neural level (e.g. modelling the spike response curves of neurons to stimulus), through the neural cluster level (e.g. modelling the release and effects of dopamine in the basal ganglia) to the complete organism (e.g. behavioral modeling of the organism's response to stimuli).

Applications:

The utility of artificial neural network models lies in the fact that they can be used to infer a function from observations and also to use it. This is particularly useful in applications where the complexity of the data or task makes the design of such a function by hand impractical.

Real life applications

The tasks to which artificial neural networks are applied tend to fall within the following broad categories: Function approximation, or regression analysis, including time series prediction and modelling. Classification, including pattern and sequence recognition, novelty detection and sequential decision making.

Data processing, including filtering, clustering, blind signal separation and compression. Application areas include system identification and control (vehicle control, process control), game-playing and decision making (backgammon, chess, racing), pattern recognition (radar systems, face identification, object recognition, etc.), sequence recognition (gesture, speech, handwritten text recognition), medical diagnosis, financial applications, data mining (or knowledge discovery in databases, "KDD"), visualization and e-mail spam filtering.

Clustering:

"The process of organizing objects into groups whose members are similar in some way". Clustering is a data mining (machine learning) technique used to place data elements into

related groups without advance knowledge of the group definitions. Popular clustering techniques include k-means clustering and expectation maximization (EM) clustering.

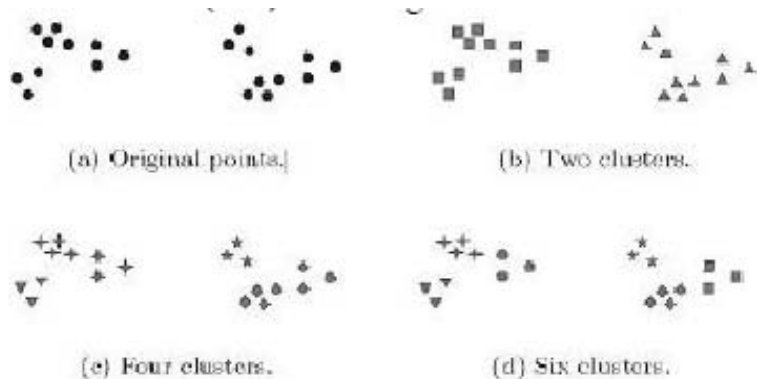


Figure 4: clustering

In this case we easily identify the 4 clusters into which the data can be divided; the similarity criterion is distance: two or more objects belong to the same cluster if they are “close” according to a given distance (in this case geometrical distance). This is called distance-based clustering. Clustering is the assignment of objects into groups (called clusters) so that objects from the same cluster are more similar to each other than objects from different clusters. Set of like elements. Elements from different clusters are not alike. The distance between points in a cluster is less than the distance between a point in the cluster and any point outside it. Problems occurring in clustering are:

- Outline handling is difficult, the elements do not naturally lie into any cluster.
- Dynamic data in the database means that cluster membership may change over time.
- Interpreting the semantic meaning of each cluster can be difficult.
- There is not a single answer to a clustering problem.

Another issue is what data should be used for clustering. We can then summarize some basic features of clustering:

- The number of clusters is not known.
- There may not be any prior knowledge concerning the clusters.
- Cluster results are dynamic.

Different Types of clusters

Clustering aims to find useful groups of objects (clusters), where usefulness is defined by the goal of the data analysis. Not surprisingly, there are several different notions of a cluster that prove useful in practice.

Well-separated A cluster is a set of objects in which each object is closer (or more similar) to every other objects in the cluster than to any object not in the cluster. Sometimes a threshold is used to specify that all the objects in a cluster must be sufficiently close (or similar) to one another.

this idealistic definition of a cluster is satisfied only when the data contains natural clusters that are quite far from each other. The distance between any two points within a group, well-separated clusters do not need to be globular, but can be any shape



Figure 5: well-separated clusters

prototype-Based A cluster is a set of objects in which each object is closer (more similar) to the prototype that defines the cluster than to the prototype of any other cluster. For data with continuous attributes, the prototype of a cluster is often a centroid, i. e., the average (mean) of all the points in the cluster. When a centroid is not meaningful, such as when the data has categorical attributes, the prototype is often named a *doid*, i. e., the most representative point of a cluster. For many types of data, the prototype can be regarded as the most central point, and in such instances, we commonly refer to prototype-based clusters as center-based clusters.

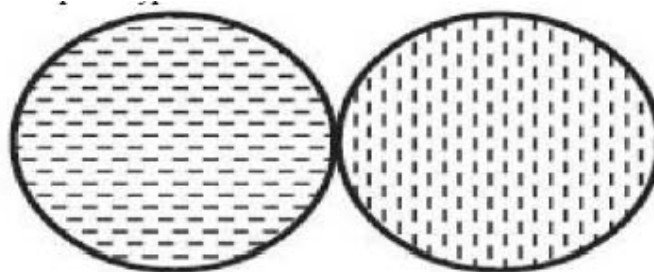


Figure 6: prototype based clusters

Graph-based if the data is represented as a graph, where the nodes are objects and the links represent connections among objects then a cluster can be defined as a connected component, i. e., a graph of objects that are connected to one another, but that have no connection to objects outside the group. An important example of graph-based clusters are contiguity-based clusters, where two objects are connected only if they are within a specified distance of each other. This implies that each object in a contiguity-based cluster is closer to some other objects in the cluster than to any point in a different cluster.

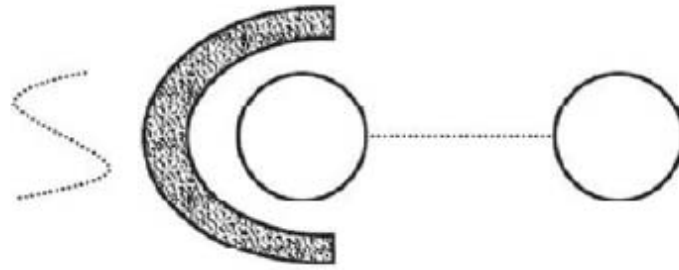


Figure 7: Graph-based clusters

Density-based clusters are dense regions of objects that are surrounded by a region of low density. A density-based definition of a cluster is often employed when the clusters are irregular or intertwined, and when noise and outliers are present.

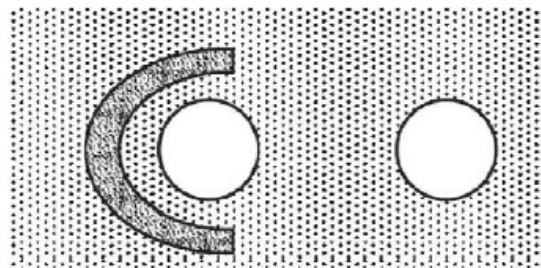


Figure 8: Density-based clusters

Shared-property (conceptual clusters) More generally, we can define a cluster as a set of objects that share some property. This definition encompasses all the previous definitions of a cluster, e.g., objects in a center-based cluster share the property that they are all closest to the same centroid or medoid.

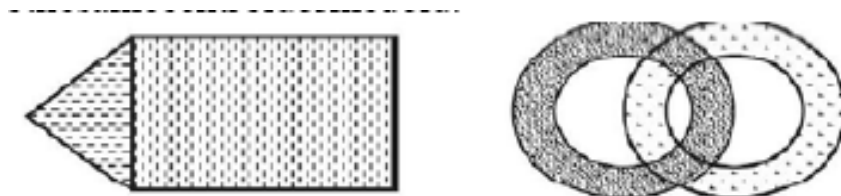


figure 9: shared-property based clusters

Noise

Clustering is often constructed on noise-free datasets. In real-world applications, it is

inevitable that the datasets contain noises, which may result in unsatisfactory results of the clustering algorithms. Outliers are sample points with values much different from those of the remaining set of data. Outliers represent error in the data or could be correct data values that are simply much different from the remaining data. A person 3.0 meter tall is exceptionally tall; this value probably would be viewed as an outlier. Some clustering techniques do not perform well with the presence of outliers. Clustering algorithms may actually search and remove outliers to ensure that they perform better. However, care must be taken in actually removing outliers. For example, suppose that the data mining problem is to predict flooding. Extremely high water level values occur infrequently, and when compared to the normal water level values may seem to be outliers. However, removing these values may not allow the data mining algorithm to work effectively because there would be no data that showed that floods ever actually occurred.

Conclusion

This paper discusses about the performance of clustering technique in the presence of noise. Noise can appear in many real world datasets and heavily corrupt the data structure. The performance of many existing algorithms is degraded by the presence of noise.

Reference

- [1] <http://sivra.in/en/datamining.html>
- [2] Varun Kumar, Nisha Rathee, " Knowledge discovery from database Using an integration of clustering and classification" (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 2, No. 3, March 2011
- [3] <http://sivra.in/en/datamining.html>
- [4] <http://vserver1.cscs.lsa.umich.edu/~spage/ONLINECOURSE/R4Decision.pdf>
- [5] Raghavendra B. K., S. K. Srivatsa, " Evaluation of logistic Regression and Neural Network Model With Sensitivity Analysis on Medical Datasets" International Journal of Computer Science and Security (IJCSS), Volume (5): Issue (5): 2011 International Journal of Computer Science and Security (IJCSS), Volume (5): Issue (5): 2011