

CHALLENGING ISSUES OF WEBMINING TECHNIQUES

R TEJASREE, Y SAMYUKTHA REDDY, S MEGHNA, B NIHARIKA, K SWATHI
STUDENTS,

Department of CSE,

Megha Institute of Engineering and Technology, Edulabad, Ghatkesar(Md), Meddchal(Dt), TS, India

Abstract–

Data mining on large databases has been a major concern in research community, due to the difficulty of analyzing huge volume of data using only traditional OLAP tools(Online Application Processing). This sort of process implies a lot of computational power , memory and disk input or output , which can only be provided by parallel computers . This research paper also conducts a formed review of application of data mining such as Education , Banking , Insurance , Medicine , Manufacturing Engineering , Health care , Transportation , Research analysis , Sales and Marketing . This paper provides a survey of various data mining techniques .These techniques are Classification analysis , Association rule learning , Anomaly or outlier detection , Clustering analysis and Regression analysis . This paper discuss the topic based on past research paper and also studies the data mining techniques and application.

INTRODUCTION

In the real world , huge amount of data are available in education , medical , industry and in many other areas .

What is Data?



- Data are raw ingredients from which statistics are created.
- Statistical analysis can be performed on data to show relationships among the variables collected.
- Through secondary data analysis, many different researchers can re-use the same data set for different purposes.

Data Mining is a type of sorting technique

which is actually used today by companies with a strong consumer focus retail , financial, communication and marketing organizations . It enables these companies to determine relationships among “Internal” factors such as price , product positioning , or staff skills , and “External” factors such as economic indicators, competitions, and customer demographics.

For example , Blockbuster entertainment mines its video rental history database to recommend rentals to individual customers. American Express can suggest products to its card holders based on analysis of their monthly expenditures.

Walmart is pioneering mass data mining to transform its supplier relationships . Walmart allows more than 3,500 suppliers, to access data on their products and perform data analysis. These suppliers use this data to identify customer buying

patterns at the store display level. In 1995, Walmart computers processed over 1 million complex data queries.

The National Basketball Association (NBA) is exploring a data mining application that can be used in conjunction with image recordings of basketball games. The advanced scout software analyzes the movements of players to help coaches orchestrate plays and strategies.

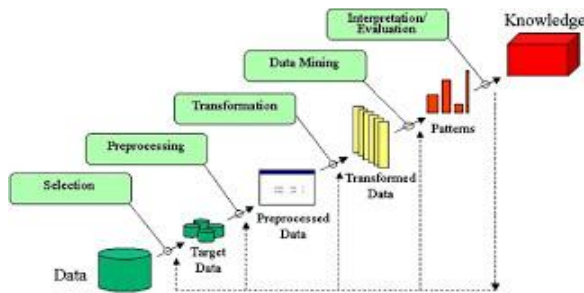


Fig: data mining database.

Data mining software analyzes relationships and patterns in stored transactions data based on open-ended user queries. Generally any four types of relationships are sought:

Classes: Stored data is used to locate data in predetermined groups.

Clusters: Data items are grouped according to logical relationship or consumer preferences.

Associations: Data can be mined to identify associations.

Sequential patterns: Data is mined to anticipate behavior patterns and trends.

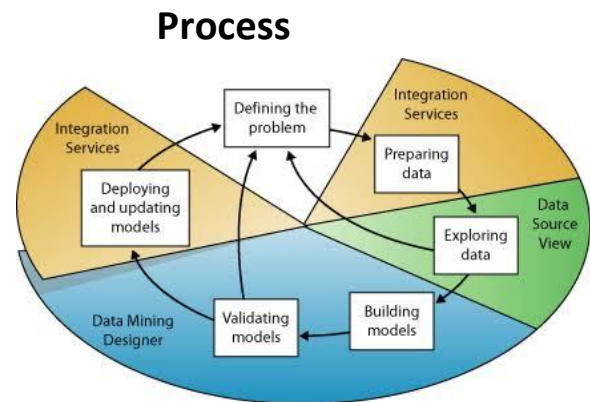


Fig: process for data mining.

DATA MINING APPLICATIONS

Various field adapted data mining technologies because of fast access of data and valuable information from a large amount of data. Data mining applications areas includes marketing, telecommunications, finance, education, medical and so on. Some of the main applications are listed below:

1. Data mining in education sector: We are applying data mining in education sector then new emerging field called "EDUCATION DATA MINING". The goals of EDM are identified as predicting students future learning behavior, studying the effects of educational support, and advancing scientific knowledge about learning. Data mining can be used by an institution to take accurate decisions and also to predict the result of the students.

2. Data mining in banking and finance: Credit and spending by customer groups can be identified by using data mining. The hidden correlations between different financial indicators can be discovered by using data mining. Data mining is used to identify customer loyalty by analyzing the data of customer purchasing activities. In the financial markets, data mining techniques such as neural

networks used in stock forecasting, price prediction and so on.

3. Data mining in insurance: Data mining is applied in claims analysis such as identifying which medical procedures are claimed together. Data mining enables forecasts which customers will potentially purchase new policies. It helps to detect fraudulent behavior. It also allows insurance companies to detect risky customers behavior patterns.

4. Data mining in medicine: Data mining enables to characterize patient activities to see incoming office visits. It helps in identifying the patterns of successful medical therapies for different illnesses. Example: smart health prediction in data mining.

5. Data mining in manufacturing engineering: Data mining tools can be very useful to discover patterns in complex manufacturing process. Data mining can be used in system-level designing to extract the relationships between product architecture, product portfolio, and customer needs data. It can also be used to predict the product development span time, cost, and dependencies among other tasks.

6. Data mining in health care: Data mining applications can be developed to evaluate the effectiveness of medical treatments. By comparing and contrasting causes, symptoms and courses of treatments, data mining can deliver an analysis of which courses of actions prove effective. In 1999, Florida hospital launched the clinical best practices initiatives with the goal of developing a standard path of care across all campuses, clinicians and patient admissions.

7. Data mining in transportations: Data mining helps determine the distribution schedules among warehouses and outlets and analyses loading patterns.

8. Data mining in research analysis: History shows that we have witnessed revolutionary changes in research. Data mining is helpful in data cleaning, data pre-processing and integration of databases. The researchers can find any similar data from the data base that might bring any change in the research. Data visualization and visual data mining provides us with a clear view of data.



Fig : data mining applications.

9. Data mining in sales and marketing: Data mining is used for “Market Basket Analysis” to provide information on what product combinations were purchased together when they were bought and in what sequence this information helps business promote their most profitable products and maximize the profit in addition to it encourages the customers to purchase related product that they may have been missed or overlooked.

TECHNIQUES

There are five types of techniques. They are

1. Classification Analysis: This analysis is used to retrieve important and relevant information about data, and metadata. It is used to classify different data in different classes. Classification is similar to clustering in a way that it also segments data records into different segments called classes. But unlike clustering, here the data analysts would have acknowledged of different classes or

cluster .So ,in classifications analysis you would apply algorithms to decide how new data should be classified . A classic example of classification analysis would be our outlook email.

2.Association Rule Learning :It refers to the method that can help you identify some interesting relations between different variables in large databases . This technique can help you unpack some hidden patterns in the data that can be used to identify variables within the data and the concurrence of different variables that appear very frequently in the dataset . Association rules are useful for examining and forecasting customer behavior.It is highly recommended in the retail industry analysis . In IT ,programmers use association rules to build programs capable of machine learning .

3.Anomaly or Outlier detection:This refers to the observation for data items in a dataset that do not match an expected pattern or an expected behavior . Anomalies are also known as outliers, novelties , noise,deviations,and exceptions . Often they provide critical and actionable information . These types of items are statistically aloof as compared to the rest of the data and hence , it indicates that something out of the ordinary has happened and requires additional attention . This technique can be used in a variety of domains ,such as intrusion detection , system health monitoring , fraud detection ,fault detection ,event detection in sensor networks ,and detecting eco-system disturbances . Analysts often remove the anomalous data from the dataset to discover results with an increased accuracy .

4. Clustering analysis:The cluster is actually a collection of data objects ; those objects are similar within the same cluster . That means the objects are similar to one another within the same group and they are rather different or they are dissimilar

or unrelated to the objects in other groups or in other clusters . Clustering analysis is the process of discovering groups and clusters in the data in such a way that the degree of association between two objects is highest if they belong to the same group and lowest otherwise . A result of this analysis can be used to create customer profiling .

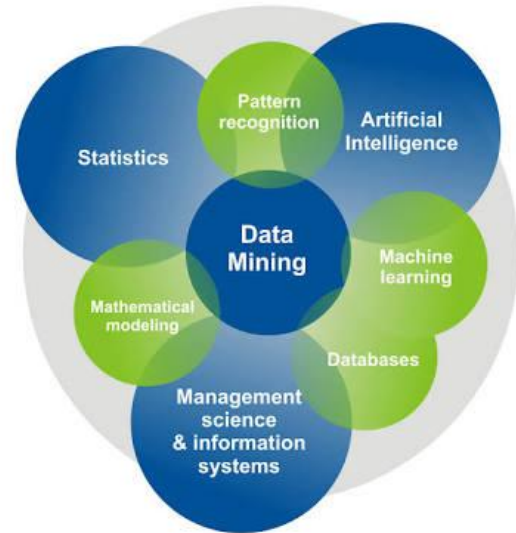


fig: techniques involved in data mining .

5.Regression analysis: In statistical terms , a regression analysis is the process of identifying and analyzing the relationship among variables . It can help you understand the characteristic value of the changes , if any one of the independent variables is varied . This means one variable is dependent on another ,but it is not vice versa . It is generally used for prediction and forecasting .

All these techniques can help analyze different data from different perspectives .

FUTURE DIRECTIONS

First , the most focused and extensively studied topic in frequent pattern mining is perhaps scalable mining methods . When we are working with data streams still it is a research challenge to derive a compact but high quality set of

patterns that are most useful in applications. The set of frequent patterns derived by most of the current patterns mining methods including ours give approximate patterns as stream is flowing continuously and some data is lost in the process of analyzing the stream.

To make frequent patterns mining an essential task in data mining, much research is needed to further develop pattern-based mining methods. For example, classification is an essential task in data mining. Construction of better classification models using frequent patterns than most other classification methods is again a research issue.

Another major research area in frequent mining is interpretation of patterns i.e., semantic annotation for frequent patterns, and contextual analysis of frequent patterns. The semantics of a frequent patterns includes deeper information. What is the meaning of the patterns; What are the synonym patterns; and What are the typical transactions that this patterns resides?

On one side, it is important to go to the core part of patterns mining algorithms, And analyze the theoretical properties of different solutions. Much work is needed to explore new applications of frequent patterns mining. For example, bioinformatics has raised a lot of challenging problems, and we believe frequent pattern mining may contribute a good deal to it with further research efforts.

Achievements: In this these our objective was to:

- Construct synopsis of data stream of transactions.
- Mine frequent itemsets.
- Mine frequent patterns.
- Mine infrequent patterns.

Construct synopsis of data stream of transactions :

The different techniques related to synopsis construction with special emphasis on reservoir sampling. We have proposed two algorithms based on reservoir sampling to construct synopsis and to mine frequent itemsets.

Mine frequent itemsets: We have proposed a new counter based algorithm to mine frequent itemsets. This work is published in an international journal.

Mine frequent patterns : In this we proposed a new data structure called Dynamic FP-tree to mine frequent patterns. Experiments have the efficiency of dynamic fp-tree.

Mine infrequent patterns: we have proposed a new algorithm based on dynamic FP-tree to mine infrequent patterns. This work is published in an international journal.

The silent features of the research work carried out are:

- A thorough literature review has been carried out.
- Data stream is continuous flow of data. Since it is not possible to analyze the whole stream. Hence a new algorithm based on reservoir sampling has been proposed to construct synopsis of data stream.
- Frequent itemset mining has number of scientific and commercial applications. This algorithm is also based on reservoir sampling.

Conclusion

Frequent pattern mining has been a focused theme in data mining research for over a decade. Abundant literature has been dedicated to this research and tremendous progress has been made

, ranging from efficient and scalable algorithms for frequent itemset mining in transaction databases to numerous research frontiers, such as sequential pattern mining, structured pattern mining, correlation mining, associative classification, and frequent pattern-based clustering, as well as their broad applications. It is believed that frequent pattern mining research has substantially broadened the scope of data analysis and will have deep impact on data mining methodologies and applications in the long run. However, there are still some challenging research issues that need to be solved before frequent pattern mining can claim a cornerstone approach in data mining applications.

REFERENCES

1. Hand D., Mannila H. and Smyth P., Principles of Data Mining, MIT Press, 2001.
2. Achlioptas, P., Scholkopf, B., and Borgwardt, K. (2011). In ACM SIGKDD Conference on data mining.
3. Azencott, C., Grimm, D., Sugiyama, M., Kawahara, Y., and Borgwardt, K. M. (2013). Bioinformatics.
4. Bishop, C. M. (2006). Pattern Recognition