

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/2383851>

A Review of Data Mining Techniques

Article in *Industrial Management & Data Systems* · October 2001

DOI: 10.1108/02635570110365989 · Source: CiteSeer

CITATIONS

188

READS

3,897

2 authors:



Sang-Jun Lee

National Fisheries Research and Development Institution

453 PUBLICATIONS 10,873 CITATIONS

[SEE PROFILE](#)



Keng Siau

City University of Hong Kong

550 PUBLICATIONS 11,517 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Supply Chain Management [View project](#)



InGaAsSb based photodetectors and Light emitting diodes for gas sensor, and emitter applications [View project](#)

A review of data mining techniques

Sang Jun Lee

University of Nebraska-Lincoln, Lincoln, Nebraska, USA

Keng Siau

University of Nebraska-Lincoln, Lincoln, Nebraska, USA

Keywords

Data mining,
Artificial intelligence, Algorithms,
Decision trees

Abstract

Terabytes of data are generated everyday in many organizations. To extract hidden predictive information from large volumes of data, data mining (DM) techniques are needed. Organizations are starting to realize the importance of data mining in their strategic planning and successful application of DM techniques can be an enormous payoff for the organizations. This paper discusses the requirements and challenges of DM, and describes major DM techniques such as statistics, artificial intelligence, decision tree approach, genetic algorithm, and visualization.

Introduction

The technologies for generating and collecting data have been advancing rapidly. At the current stage, lack of data is no longer a problem; the inability to generate useful information from data is! The explosive growth in data and database results in the need to develop new technologies and tools to process data into useful information and knowledge intelligently and automatically. Data mining (DM), therefore, has become a research area with increasing importance (Weiss and Indurkha, 1998; *Technology Forecast*, 1997; Fayyad *et al.*, 1996; Piatetsky-Shapiro and Frawley, 1991).

DM is the search for valuable information in large volumes of data (Weiss and Indurkha, 1998). It is the process of nontrivial extraction of implicit, previously unknown and potentially useful information such as knowledge rules, constraints, and regularities from data stored in repositories using pattern recognition technologies as well as statistical and mathematical techniques (*Technology Forecast*, 1997; Piatetsky-Shapiro and Frawley, 1991). Many companies have recognized DM as an important technique that will have an impact on the performance of the companies.

DM is an active research area and research is ongoing to bring statistical analysis and artificial intelligence (AI) techniques together to address the issues.

Current trends on data mining

Just five years ago, only 50 researchers took part in the knowledge discovery and data mining conference workshop. Today, however, knowledge discovery nuggets, the well-known monthly electronic newsletter by Gregory Piatetsky-Shapiro, has more than

4,000 readers. Moreover, data mining continues to attract more and more attention in the business and scientific communities. In a 1997 report, Stamford, Connecticut-based Gartner Group mentioned: "Data mining and artificial intelligence are at the top of the five key technology areas that will clearly have a major impact across a wide range of industries within the next three to five years." Many companies currently use computers to capture details of business transactions such as banking and credit card records, retail sales, manufacturing warranty, telecommunications, and myriad other transactions. Data mining tools are then used to uncover useful patterns and relationships from the data captured.

Currently, data mining techniques, tools, and researches are being expanded to the various fields. For example, the DM tool, intelligent text-mining system, extracts text fragments relevant to user queries, automatically creates and processes a series of new queries, and assembles a new text. The output enables the user to see the new aspects of a given theme. This tool is a rule-based system using complex heuristics.

Data warehousing is one of the most important research areas related to DM. A data warehouse is a read-only database developed for analyzing business situations and supporting decision makers. The data warehouse includes large volumes of subject-oriented data, where all levels of an organization can find the information in a timely manner. DM goes together with the data warehousing which is necessary to organize historical information gathered from large-scale client/server-based applications. In other words, DM can add values to the information assets of organizations in different sectors, through effective induction of large corporate data warehouses into a client-server. Therefore, developing an advanced client-server induction system capable of supporting efficient and effective data mining of large



databases in business environment is one of the active research areas.

Requirements and challenges of DM

DM is a relatively new field and there are many challenges to be faced. Extracting useful information from data can be a complicated and sometimes a difficult process. In this section, we look at some of the requirements and challenges of data mining (adapted from Chen *et al.*, 1996).

Ability to handle different types of data

Many database systems have complex data types, such as hypertext, multimedia data, and spatial data. If a DM technique is robust and powerful, it should be able to perform effective DM on various types of data structures. Though ideal, it is impractical to expect a DM technique to handle all kinds of data and to perform different goals of DM effectively. In general, a specific DM system is built for mining knowledge from a specific kind of data.

Graceful degeneration of DM algorithms

The DM algorithms should be efficient and scaleable. The performance of the algorithm should degenerate gracefully. In other words, the searching, mining, or analyzing time of a DM algorithm should be predictable and acceptable as the size of the database increases.

Valuable DM results

DM system should be able to handle noise and exceptional data efficiently. The discovered information must precisely depict the contents of the database and be beneficial for certain applications. Also, the quality of the discovered information should be interesting and reliable.

Representation of DM requests and results

DM identifies facts or conclusions based on sifting through the data to discover patterns or anomalies (*Technology Forecast*, 1997). To be effective, the systems should allow users to discover information from their own perspectives and the information should be presented to the users in forms that are comfortable and easy to understand. High-level query languages or graphical user interface is required to express the DM requests and the discovered information. End users should be able to specify task commands for the DM system and the results from the DM system should be understandable and usable.

Mining at different abstraction levels

It is very difficult to specify exactly what to look for in a database or how to extract useful information from a database. Besides, the value of a piece of information is in the eyes of the beholder – one person's "gold mine" could easily be another person's garbage. To facilitate the mining process, the systems should allow the users to mine at different abstraction levels. For example, a high-level query might disclose an interesting trace that warrants further exploration. Thus, it is important for DM tools to support mining at different levels of granularity.

Mining information from different sources of data

In the ages of the Internet, Intranets, Extranets, and data warehouses, many different sources of data in different formats are available. Mining information from heterogeneous database and new data formats can be challenges in DM. The DM algorithms should be flexible enough to handle data from different sources.

Protection of privacy and data security

DM is a threat to privacy and data security because when data can be viewed from many different angles at different abstraction levels, it threatens the goal of keeping data secured and guarding against the intrusion on privacy. For example, it is relatively easy to compose a profile of an individual (e.g. personality, interests, spending habits, etc.) with data from various sources.

Data mining steps

In general, there are three main steps in DM: preparing the data, reducing the data and, finally, looking for valuable information. The specific approaches, however, differ from companies to companies and researchers to researchers. For example, IBM (reported in *Technology Forecast*, 1997) defined four major operations for DM:

- 1 *Predictive modeling*: using inductive reasoning techniques such as neural networks and inductive reasoning algorithms to create predictive models.
- 2 *Database segmentation*: using statistical clustering techniques to partition data into clusters.
- 3 *Link analysis*: identifying useful associations between data.
- 4 *Deviation detection*: detecting and explaining why certain records cannot be put into specific segments.

Fayyad *et al.* (1996), on the other hand, proposed the following steps:

- 1 Retrieving the data from a large database.
- 2 Selecting the relevant subset to work with.
- 3 Deciding on the appropriate sampling system, cleaning the data and dealing with missing fields and records.
- 4 Applying the appropriate transformations, dimensionality reduction, and projections.
- 5 Fitting models to the preprocessed data.

Classifying DM techniques

Many DM techniques and systems have been developed and designed. These techniques can be classified based on the database, the knowledge to be discovered, and the techniques to be utilized. In this section, we review one of the classification schemes proposed by Chen *et al.* (1996).

Based on the database

There are many database systems that are used in organizations, such as relational database, transaction database, object-oriented database, spatial database, multimedia database, legacy database, and Web database. A DM system can be classified based on the type of database it is designed for. For example, it is a relational DM system if the system discovers knowledge from relational database and it is an object-oriented DM system if the system finds knowledge from object-oriented database.

Based on the knowledge

DM systems can discover various types of knowledge, including association rules, characteristic rules, classification rules, clustering, evolution, and deviation analysis. DM systems can also be classified according to the abstraction level of the discovered knowledge. The knowledge may be classified into general knowledge, primitive-level knowledge, and multiple-level knowledge.

Based on the techniques

DM systems can also be categorized by DM techniques. For example, a DM system can be categorized according to the driven method, such as autonomous knowledge mining, data-driven mining, query-driven mining, and interactive DM techniques. Alternatively, it can be classified according to its underlying mining approach, such as generalization-based mining, pattern-based mining, statistical- or mathematical-based mining and integrated approaches.

Major DM techniques

In this section, we review and discuss the major DM techniques.

Statistics

Statistics is an indispensable component in data selection, sampling, DM, and extracted knowledge evaluation. It is used to evaluate the results of DM to separate the good from the bad. In data cleaning process, statistics offer the techniques to detect “outliers”, to smooth data when necessary, and to estimate noise. Statistics can also deal with missing data using estimation techniques.

Techniques in clustering and designing of experiments come into play for exploratory data analysis. Work in statistics, however, has emphasized generally on theoretical aspects of techniques and models. As a result, search, which is crucial in DM, has received little attention. In addition, interface to database, techniques to deal with massive data sets, and techniques for efficient data management are very important issues in DM. These issues, however, have only begun to receive attention in statistics (Kettenring and Pregibon, 1996).

Techniques for mining transactional/relational database

Mining association rules in transactional or relational database has been the most attractive topics in database field (Agrawal *et al.*, 1993; Han and Fu, 1995; Mannila *et al.*, 1994; Savasere *et al.*, 1995; Srikant and Agrawal, 1995). The task is to derive a set of strong association rules in the form of “ $A_1 \wedge \dots \wedge A_m \Rightarrow B_1 \wedge \dots \wedge B_n$,” where A_i (for $i \in \{1, \dots, m\}$) and B_j (for $j \in \{1, \dots, n\}$) are attribute-value sets, from the associated data sets in a database. For example, one may find the association rule: if a customer buys one brand of beer, he/she usually buys another brand of chips in the same transaction. Because mining association rules might require scanning through a massive transaction database repeatedly, the required processing power could be enormous.

Artificial intelligence (AI) techniques

AI techniques are widely used in DM. Techniques such as pattern recognition, machine learning, and neural networks have received much attention. Other techniques in AI such as knowledge acquisition, knowledge representation, and search, are relevant to the various process steps in DM.

Classification is one of the major DM problems. Classification is the process of dividing a data set into mutually exclusive groups such that the members of each group

are as “close” as possible to one another, and the members of different groups are as “far” as possible from one another. For example, a typical classification problem is to divide a database of customers into groups that are as homogeneous as possible with respect to a variable such as creditworthiness. One solution to the classification problem is to use neural network. According to Lu *et al.* (1996), neural network-based DM approach consists of three major phases:

- 1 *Network construction and training*: in this phase, a layered neural network based on the number of attributes, number of classes, and chosen input coding method are trained and constructed.
- 2 *Network pruning*: in this phase, redundant links and units are removed without increasing the classification error rate of the network.
- 3 *Rule extraction*: classification rules are extracted in this phase.

Other AI techniques that can be used for DM include case-based reasoning and intelligent agents. Case-based reasoning uses historical cases to recognize patterns and the intelligent agent approach employs a computer program (i.e. an agent) to sift through data.

Decision tree approach

Decision trees are tree-shaped structures that represent sets of decisions. The decision tree approach can generate rules for the classification of a data set. Specific decision tree methods include Classification and Regression Trees (CART) and Chi Square Automatic Interaction Detection (CHAID). CART and CHAID are decision tree techniques used for classification of a data set. They provide a set of rules that can be applied to a new (unclassified) data set to predict which records will have a given outcome. CART typically requires less data preparation than CHAID.

Genetic algorithm

Genetic algorithm is a relatively new software paradigm inspired by Darwin’s theory of evolution. A population of rules, each representing a possible solution to a problem, is initially created at random. Then pairs of rules (usually the strongest rules are selected as parents) are combined to produce offspring for the next generation. A mutation process is used to randomly modify the genetic structures of some members of each new generation. The system runs for dozens or hundreds of generations. The process is terminated when an acceptable or optimum solution is found, or after some fixed time limit. Genetic algorithms are appropriate for

problems that require optimization with respect to some computable criterion. This paradigm can be applied to DM problems. The quantity to be minimized is often the number of classification errors on a training set. Large and complex problems require a fast computer in order to obtain appropriate solutions in a reasonable amount of time. Mining large data sets by genetic algorithms has become practical only recently due to the availability of affordable high-speed computers.

Visualization

A picture is worth thousands of numbers! Visual DM techniques have proven the value in exploratory data analysis, and they also have a good potential for mining large database. This approach requires the integration of human in the DM process. Tufte (1983, 1990) provided many examples of visualization techniques that have been extended to work on large data sets and produce interactive displays. There are several well-known techniques for visualizing multidimensional data sets: scatterplot matrices, coplots, projection matrices, parallel coordinates, projection matrices, and other geometric projection techniques such as icon-based techniques, hierarchical techniques, graph-based techniques, and dynamic techniques.

Managerial implications, directions and strategies

The real question is what DM can do in real business field? DM can give solutions for banking and credit industry such as credit scoring, fraud detection, and customer segmentation. In today’s competitive environment, maintaining widely-used bureau based tools may not be an effective long-term strategy. However, DM can develop proprietary bureau scores for general risk, bankruptcy, revenue and response, among others, that are not available to competitors and are optimized for the business needs. Tools adapting DM techniques provide credit managers and underwriters for rapidly building and analyzing credit application scorecards allowing increased acceptances and reduced bad debt. When provided with a history of fraudulent and genuine credit applications, DM tools learn the patterns of fraud and subsequently identify fraud successfully.

Nowadays, the Internet is one of the most important markets and channels in the business environment. It is so important to deliver the right message at the right time to

the right people. In fact, the Internet environment, especially, World Wide Web (WWW) has very unstructured data format from a DM point of view. DM has been applied to the structured databases. However, the potential of Web mining to help people navigate, search, and visualize the contents of Web is enormous, and the Web mining is feasible in practice (Etzioni, 1996). DM tools adapting the Internet technology can provide a targeted banner campaign, rich media, a storefront, a promotion, or direct e-mail. Different services are provided to different people. DM tools anticipate customers' needs before, during, and after running a campaign. It reports behavioral data about Web visitors, such as who is visiting, what content they access, where they came from, and why they purchase. The knowledge gained from the tool can be used to measure return on investment (ROI) of marketing campaigns and make better online business decisions. One second can turn a browser into a buyer. The key is making the right suggestion to the right buyer at the right time. It is called suggestive selling, and DM tools can do the job in this Internet arena.

DM has been recognized as one of the most important techniques in e-commerce market – especially in providing filtering techniques that are key to success of today's top Internet marketers. They drive the personalized recommendations that turn more site browsers into buyers, increase cross-selling and up-selling, and deepen customers' loyalty with every purchase. In e-commerce marketing, most of the market analysis tools rely on data mining and data warehousing. These technologies are used to uncover relationships among products – “product affinities” and between customers and products – “market segmentation”.

DM techniques are even used for sports and entertainment solutions. IBM has prototyped leading DM technology to help National Basketball Association (NBA) coaches and league officials organize and interpret the mountains of data amassed at every game. Using DM software called Advanced Scout to prepare for a game, a coach can quickly review countless statistics: shots attempted, shots blocked, assists made, and personal fouls. The tool can also detect patterns in these statistics that a coach may not know about. For example, the tool may uncover information such as this player is most effective with these players and under these circumstances.

Nowadays, the use of DM technology is widespread in any industry. A large multinational retailer uses DM technique to refine inventory stocking levels, by store and by item, to dramatically reduce out-of-stock or

overstocking situations and thereby improve revenues and reduce forced markdowns. A health maintenance group uses DM technique to predict which of its members are most at risk of specific major illnesses. This presents opportunities for timely medical intervention and preventative treatment to promote the patients' well-being and reduce the healthcare provider's costs.

Conclusion

Having the right information at the right time is crucial for making the right decision. The problem of collecting data, which used to be a major concern for most organizations, is almost resolved. In the millennium, organizations will be competing in generating information from data and not in collecting data. Industry surveys indicated that over 80 percent of *Fortune 500* companies believe that data mining would be a critical factor for business success by the year 2000 (Baker and Baker, 1998). Obviously, DM will be one of the main competitive focuses of organizations. Although progresses are continuously been made in the DM field, many issues remain to be resolved and much research has to be done.

References

- Agrawal, R., Imielinski, T. and Swami, A. (1993). *Mining Association Rules between Sets of Items in Large Databases*, Paper presented at the ACM SIGMOD, May.
- Baker, S. and Baker, K. (1998), “Mine over matter”, *Journal of Business Strategy*, Vol. 19 No. 4, pp. 22-7.
- Chen, M.S., Han, J. and Yu, P. (1996), “Data mining: an overview from a database perspective”, *IEEE Transactions on Knowledge and Data Engineering*, Vol. 8 No. 6, pp. 866-83.
- Etzioni, O. (1996), “The World-Wide Web: quagmire or gold mine?”, *Communication of the ACM*, Vol. 39 No. 11, pp. 65-8.
- Fayyad, U., Djorgovski, S.G. and Weir, N. (1996), “Automating the analysis and cataloging of sky surveys”, in Fayyad, U., Piatetsky-Shapiro, G., Smyth, P. and Uthurusamy, R. (Eds), *Advances in Knowledge Discovery and Data Mining*, MIT Press, Cambridge, MA, pp. 471-94.
- Fayyad, U., Piatetsky-Shapiro, G. and Smyth, P. (1996), “From data mining to knowledge discovery: an overview”, in Fayyad, U., Piatetsky-Shapiro, G., Smyth, P. and Uthurusamy, R. (Eds), *Advances in Knowledge Discovery and Data Mining*, MIT Press, Cambridge, MA.
- Han, J. and Fu, Y. (1995), *Discovery of Multiple-Level Association Rules from Large Databases*, Paper presented at the 21st Int'l Conf. Very Large Data Bases, September.

- Kettenring, J. and Pregibon, D. (1996), *Committee on Applied and Theoretical Statistics: Workshop on Massive Data Sets*, Paper presented at the National Research Council, Washington, D.C.
- Lu, H., Setiono, R. and Liu, H. (1996), “Effective Data Mining Using Neural Networks”, *IEEE Transactions on Knowledge and Data Engineering*, Vol. 8 No. 6, pp. 957-61.
- Mannila, H., Toivonen, H. and Verkamo, A.I. (1994), *Effective Algorithms for Discovering Association Rules*, paper presented at the AAAI Workshop, Knowledge Discovering in Databases, July.
- Piatetsky-Shapiro, G. and Frawley, W.J. (1991), *Knowledge Discovery in Database*, AAAI/MIT Press.
- Savasere, A., Omiecinski, E. and Navathe, S. (1995), *An Effective Algorithm for Mining Association Rules in Large Databases*, paper presented at the 21st International Conference, Very Large Data Bases, September.
- Srikant, R. and Agrawal, R. (1995), *Mining Generalized Association Rules*. Paper presented at the 21st International Conference, Very Large Data Bases, September.
- Technology Forecast: 1997* (1997), Price Waterhouse World Technology Center, Menlo Park, CA
- Tufte, E.R. (1983), *The Visual Display of Quantitative Information*, Graphics Press, Cheshire, CN.
- Tufte, E.R. (1990), *Envisioning Information*, Graphics Press, Cheshire, CN.
- Weiss, S.H. and Indurkha, N. (1998), *Predictive Data Mining: A Practical Guide*, Morgan Kaufmann Publishers, San Francisco, CA.