

A Review Paper on Data Mining Techniques with Duo Mining

D. Sheila Freeda

AP / MCA

Department of Master of Computer Applications

Er. Perumal Manimekalai College of Engineering, Hosur, Tamil Nadu, INDIA

Abstract - A very large amount of data stored in databases is increasing at a tremendous speed. This requires a need for new techniques and tools to aid humans in automatically and intelligently analyzing large data sets to acquire useful information. This growing need gives a view for a new research field called Knowledge Discovery in Databases (KDD) or Data Mining, which attract a attention from researchers in many different fields including database design, statistics, pattern recognition, machine learning, and data visualization. The combination of data and text mining is referred to as “Duo-mining”. Text and data mining are fast growing areas and are believed to have high commercial potential value in knowledge discovery and information filtering areas of application. Although text mining manages unstructured data, most of knowledge discovery and information filtering can be done using data mining.

Keywords— Data Mining, classification, Duo Mining, Pattern Recognition, Text Mining

I. INTRODUCTION

The last decade has experienced a revolution in information availability and exchange of it through internet. In the same strength more business as well as organizations began to collect data related to their own operations, while the database technologist have been seeking efficient mean of storing, retrieving and manipulating data, the machine learning community focused on techniques which used for developing, learning and acquiring knowledge from the data. Data Mining is the process of analyzing data from different perspectives and summarizing it into useful information.

Data mining consists of extract, transform, and load transaction data onto the data warehouse system, store and manage the data in a multidimensional database system, by using application software analyses the data, provide data access to business analysts and information technology professionals, present the data in a useful format, like a graph or table. Data mining involves the anomaly detection, association, classification, regression, rule learning, summarization and clustering.

II. DATA MINING

Data mining is the exploration and analysis of large data sets, in order to discover meaningful pattern and rules. The key idea is to find effective way to combine the computer’s power to process the data with the human eye’s ability to detect patterns. The objective of data mining is to design and work efficiently with large data sets. Data mining is the databases. In data mining the data can be mined by component of wider process called knowledge discovery from database. Data Mining is the process of analyzing data from different perspectives and summarizing the results as useful information. It has been defined as "the nontrivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data"

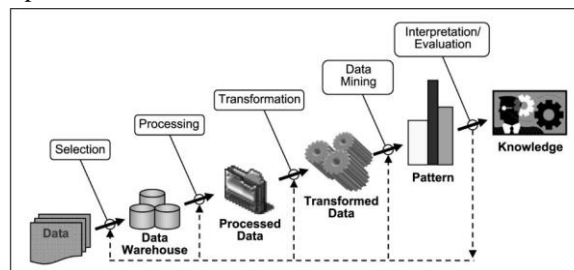


Fig. 1 : Steps in Data Mining process

The definition of data mining is closely related to another commonly used term knowledge discovery . Data mining is an interdisciplinary, integrated database, artificial intelligence, machine learning, statistics, etc. Many areas of theory and technology in current era are databases, artificial intelligence, data mining and statistics is a study of three strong large technology pillars. Data mining is a multi-step process, requires accessing and preparing data for a mining the data, data mining algorithm, analyzing results and taking appropriate action. The data, which is accessed can be stored in one or more operational databases. In data mining the data can be mined by various processing.

2.1. Supervised Learning

In supervised learning (often also called directed data mining) the variables under investigation can be split into two groups: explanatory variables and one (or more) dependent variables. The goal of the analysis is to specify a relationship between the dependent variable and explanatory variables the as it is done in regression analysis. To proceed with directed data mining techniques the values of the dependent variable must be known for a sufficiently large part of the data set.

2.2. Unsupervised Learning:

In unsupervised learning, all the variables are treated in same way, there is no distinction between dependent and explanatory variables. However, in contrast to the name undirected data mining, still there is some target to achieve. This target might be as data reduction as general or more specific like clustering. The dividing line between unsupervised learning and supervised learning is the same that distinguishes discriminate analysis from cluster analysis. Supervised learning requires, target variable should be well defined and that a sufficient number of its values are given. Unsupervised learning typically either the target variable has only been recorded for too small a number of cases or the target variable is unknown.

III. TASKS OF DATA MINING

Data mining as a term used for the specific classes of six activities or tasks as follows:

1. Classification
2. Estimation
3. Prediction
4. Affinity grouping or association rules
5. Clustering
6. Description and visualization

The first three tasks - classification, estimation and prediction rules are examples of directed data mining or supervised learning. In directed data mining, the goal is to use the available data to build a model that describes one or more particular attribute(s) of interest (target attributes or class attributes) in terms of the rest of the available attributes. The next three tasks – association rules, clustering and description are examples of undirected data mining i.e. no attribute is singled out as the target, the main goal is to establish some relationship among all attributes.

3.1. Classification

Classification consists of examining the features of a newly presented object and assigning to it a predefined class. The classification task is characterized by the well-defined classes, and a training set consisting of reclassified examples. The task is to build a model that can be applied to unclassified data in order to classify it. Examples of classification tasks include:

- Classification of credit applicants as low, medium or high risk
- Classification of mushrooms as edible or poisonous
- Determination of which home telephone lines are used for internet access

3.2. Estimation

Estimation deals with continuously valued outcomes. Given some input data, we use estimation to come up with a value for some unknown continuous variables such as income, height or credit card balance. Some examples of estimation tasks include:

- Estimating the number of children in a family from the input data of mothers' education
- Estimating total household income of a family from the data of vehicles in the family
- Estimating the value of a piece of a real estate from the data on proximity of that land from a major business center of the city.

3.3. Prediction

Any prediction can be thought of as classification or estimation. The difference is one of emphasis. When data mining is used to classify a phone line as primarily used for internet access or a credit card transaction as fraudulent, we do not expect to be able to go back later to see if the classification was correct. Our classification may be correct or incorrect, but the uncertainty is due to incomplete knowledge only: out in the efforts, it is possible to check. Predictive tasks feel different because the records are classified according to some predicted future behavior or way to check the

accuracy of the classification is to wait and see. Examples of prediction tasks include:

- Predicting the size of the balance that will be transferred if a credit card prospect accepts a balance transfer offer
- Predicting which customers will leave within next six months
- Predicting which telephone subscribers will order a value-added service such as three-way calling or voice mail.

Any of the techniques used for classification and estimation can be adopted for use in prediction by using training examples where the value of the variable to be predicted is already known, along with historical data for those examples. The historical data is used to build a model that explains the current observed behavior. When this model is applied to current inputs, the result is a prediction of future behavior.

3.4. Association Rules

An association rule is a rule which implies certain association relationships among a set of objects (such as “occur together” or “one implies the other”) in a database. Given a set of transactions, where each transaction is a set of literals (called items), an association rule is an expression of the form $X \Rightarrow Y$, where X and Y are sets of items. The intuitive meaning of such a rule is that transactions of the database which contain X tend to contain Y . An example of an association rule is: “30% of farmers that grow wheat also grow pulses; 2% of all farmers grow both of these items”. Here 30% is called the confidence of the rule, and 2% the support of the rule. The problem is to find all association rules that satisfy user-specified minimum support and minimum confidence constraints.

3.5. Clustering

Cluster analysis can be used as a standalone data mining tool to gain insight into the data distribution, or as a pre-processing step for other data mining algorithms operating on the detected clusters. Many clustering algorithms have been developed and are categorized from several aspects such as partitioning methods, hierarchical methods, density-based methods, and grid-based methods. Further data set can be numeric or categorical. Clustering is the task of segmenting a diverse group into a number of similar subgroups or clusters. What distinguishes clustering from classification is that clustering does not rely on predefined classes. In clustering, there are no predefined classes. The records are grouped together on the basis of self-similarity. Clustering is often done as a prelude to some other form of data mining or modeling. For example, clustering might be the first step in a market segmentation effort, instead of trying to come up with a one-size-fits-all rule for determining what kind of promotion works best for each cluster[6].

DUO-MINING

Duo-Mining is the variation of data and text mining. It has demonstrated especially well for the banking and credit card companies in order to take better decisions. As separate capabilities, of the pattern finding technologies of data mining and text mining have been around for years. However, it is only recently that enterprises have been started to use the two in acycle - and have discovered that it is a combination that is worth more than the sum of its parts.

They are similar because they both "mine" large

amounts of data, and looking for significant patterns. However, what they evaluate is quite different. Instead of only being able to analyze the structured data they collect from transactions, they can add call logs from customer services and further analyze customers

and spending patterns from the text mining side. These new developments in text mining technology that go beyond simple searching methods are the key to information discovery which is generally work on the unstructured data.

There are several method of data mining which handle the following or application of mining:

- Spatial mining
- Multimedia mining
- Text mining
- Web mining

Spatial mining:

Spatial is basically a three-dimensional object, and mining is extraction of patterns. Non-trivial searches “robotic” as possible to diminish human effort. It refers to the extraction of knowledge, spatial relationship, or other fascinating patterns not explicitly stored in spatial databases. Such mining demands an incorporation of data mining with spatial database technology.

Multimedia mining:

Its stores and manages a large collection of multimedia data, such as audio, video, image, image, hypertext data, which contain text, text markups and linkage. Multimedia database system is gradually more common due to the trendy use of audio video equipment, digital cameras and the internet.

Text mining:

In text mining, the goal is to discover unidentified information, something that no one yet knows and so could not have yet written down. Text mining is a distinction on a field called data mining that tries to find motivating patterns from bulky databases.

Web mining:

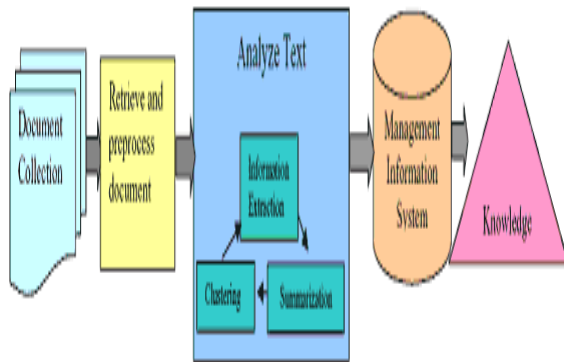
The application of data mining techniques to find out patterns from the Web. According to analysis targets, web mining can be divided into three are Web usage mining, Web content mining and Web Structure mining.

We can use data mining in so many different kinds of places in the world are as follows:

- Analysts and managers who deals with strategic and tactical decision making.
- Managers responsible for revenue and cost detection.
- Risk managers in insurance, to minimize the risk of claim and to maximize the profit.
- Educators to improve educational processes to conduct researchers, provide analysis of education effectiveness and institutional decision.
- Scientist to provide new knowledge for researchers in various fields.
- Use knowledge discovery for various opportunities like sales, forecasting, market researches etc.

IV. TEXT MINING AND ITS TYPES WITH TRADITIONAL TECHNIQUE:

It is the process of extracting fascinating and nontrivial information and knowledge from unstructured text. Text mining has been defined as “the discovery by computer of new, previously unknown information, by automatically extracting information from different written resources” Text mining is similar to data mining, except that data mining gear are designed to handle structured data from databases or XML files, but text mining can work with unstructured or semi-structured data sets such as emails, full-text documents, HTML files, etc. As a result, text mining is a much better solution for companies, where large volumes of diverse types of information must be multipart and managed.



Types of Text mining:

- Text classification.
- Text clustering.
- Keyword based association rule.

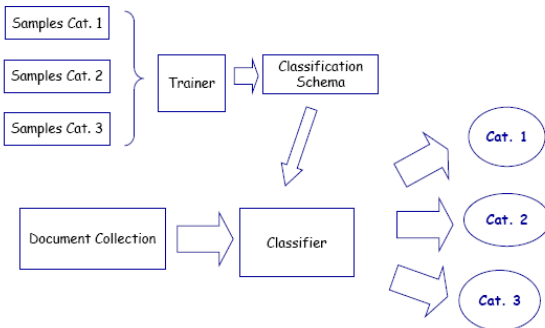
Text classification: Involuntary classification for the bulky number of online text documents (web pages, emails, corporate intranets etc.). Text document classification is differs from the classification of relational statistics. Documents databases are not prearranged according to attribute value pairs.

Following steps are taken in the process of Classification

- Data preprocessing.

- Description of training set and test sets.
- Creation of the classification model using the preferred classification algorithm.
- Classification model substantiation.
- Classification of new/unknown text documents.

Classification schema:



Text clustering: The process of isolating a dataset into reciprocally elite groups such that the members of each group are as "close" as possible to one another, and unlike groups are as "far" as possible from one another, where distance is considered with esteem to all available variables.

Keyword based association rule: Collect sets of keywords or terms that take place currently collectively and then find the association or parallel relationships among them.

Though doing keyword based association rule , we include to follow some variety of steps to analysis it:

- Preprocess the text data by parsing, stemming, removing stop words, etc.
- Evoke association mining algorithms
 - o Consider each document as a transaction.
 - o View a set of keywords in the document as a set of items in the transaction Terms level association mining
 - o No need for human effort in tagging documents.
 - o The number of meaningless result and the execution time is greatly reduced.

In this paper we have illustrated the various methods of Data mining. The proper use of Duo Mining is a combination of text and data mining have been illustrated. It also concludes the proper analysis of text and data mining in correlation with Techniques, methods, process and architecture. The use of text mining and its various types has been correlated with data mining.

REFERENCES

- [1] Donald Michie, Data Mining Discovering Interesting Relationships in Large Data Sets, Retrieved from <http://www.aaai.org/aitopics/pmwiki/pmwiki.php/AITopics/DataMining>
- [2] Wikipedia, free encyclopedia, Data mining, Retrieved from http://en.wikipedia.org/wiki/Data_mining
- [3] Discovering hidden value in your data warehouse, Retrieved from <http://www.theartling.com/text/dmwhite/dmwhite.htm>
- [4] Mining Object, Spatial, Multimedia, Text, and Web Data, Retrieved from http://www.dataminingtools.net/wiki/applications_of_data_mining.php
- [5] Wikipedia, free encyclopedia, web mining, Retrieved from http://en.wikipedia.org/wiki/Web_mining
- [6] Biomarker information extraction tool (BIET) development using natural language processing and machine, retrieved from <http://portal.acm.org/citation.cfm?id=1741927>
- [7] Mining Text And Web Data retrieved from <http://www.slideshare.net/pierluca.lanzi/machinelearning-and-data-mining-19-mining-text> .
- [8] Schwartz, A. S. And Hearst, M. A. (2003), 'A Simple Algorithm For Identifying Abbreviation Definitions In Biomedical Text', In 'Proceedings Of The 8th Pacific Symposium On Biocomputing', 3rd-7th January, Hawaii, Pp. 451-462. 29. M
- [9] Chen, H., Lally, A. M., Zhu, B., And Chau, M. (2003). "Helpfulmed: Intelligent Searching For Medical Information Over The Internet," *Journal Of The American Society For Information Science And Technology*, 54(7), 683-694, 2003. This Article Provides An Overview Of Medical Information Retrieval Techniques On The Internet, Including Web Crawling, Co-Occurrence Analysis, And Document Visualization.
- [10] Yang, Y. And Liu, X. (1999). "A Re-Examination Of Text Categorization Methods, In *Proceedings Of The 22nd Annual International ACM Conference On Research And Development In Information Retrieval (SIGIR'99)*, 1999, Pp. 42-49.
- [11] Mining Biomedical Literature Using Information Extraction ,Ronon Feldman, Yizhar Regev, Michal Finkelstein-Landau, Eyal Hurvitz & Boris Kogan Clearforest Corp, USA & Israel
- [12] Alexander Pertsemlidis; TEXT MINING THE BIOMEDICAL LITERATURE, UT Southwestern Medical Center, 5323 Harry Hines, Boulevard, Dallas, Texas 75390-8573

- [13] Type - 1 Diabetes Mellitus: Indian And Global Scene – Burden & Challenges. Diabetes Department, Voluntary Health Services, Chennai, Tamil Nadu, India.
- [14] Locating Previously Unknown Patterns In Data-Mining Results: A Dual Data- And Knowledge-Mining Method, Mir S Siadaty* And William A Knaus, Address: Department Of
- [15] Public Health Sciences, University Of Virginia School Of Medicine, Box 800717, Charlottesville, Virginia, 22908, USA Email: Mir S Siadaty* - Mirsiadaty@Virginia.Edu; William A Knaus - Wak4b@Virginia.Edu.