

An Empirical Study on Class Association Rules Mining

Ajaya Kumar Parida, Subhendu Kumar Pani

Asso.Prof., Dept.of CSE, OEC, BBSR
Asso.Prof., Dept.of CSE , OEC, BBSR

Abstract

Association rule mining is a well-known technique in data mining. It is able to reveal all interesting relationships, called associations, in a potentially large database. However, how interesting a rule is depends on the problem a user wants to solve. Existing approaches employ different parameters to guide the search for interesting rules. Class association rules which combine association rule mining and classification are therefore concerned with finding rules that accurately predict a single target (class) variable. The key strength of association rule mining is that all interesting rules are found. The number of associations present in even moderate sized databases can be, however, very large – usually too large to be applied directly for classification purposes. Therefore, any classification learner using association rules has to perform three major steps: Mining a set of potentially accurate rules, evaluating and pruning rules, and classifying future instances using the found rule set. In this work, we study association rule mining. Using a systematic approach, we generate a set of best association rules which can show higher predictive performance. We use two most popular algorithms namely Apriori and Predictive Apriori for the purpose. The dataset is drawn from a breast cancer detection-decision context available at UCI machine learning repository.

1. Introduction

Data mining is considered to be an emerging technology that has made revolutionary change in the information world. The term 'data mining' (often called as knowledge discovery) refers to the process of analysing data from different perspectives and summarizing it into useful information by means of a number of analytical tools and techniques, which in turn may be useful to increase the performance of a system.

Technically, "data mining is the process of finding correlations or patterns among dozens of fields in large relational databases". Therefore, data mining consists of major functional elements that transform data onto data warehouse, manage data in a multidimensional database, facilitates data access to information professionals or

analysts, analyse data using application tools and techniques, and meaningfully presents data to provide useful information.

According to the Gartner Group, "data mining is the process of discovering meaningful new correlation patterns and trends by sifting through large amount of data stored in repositories, using pattern recognition technologies as well as statistical and mathematical techniques"[3] . Thus use of data mining technique has to be domain specific and depends on the area of application that requires a relevant as well as high quality data.

More precisely, data mining refers to the process of analysing data in order to determine patterns and their relationships. It automates and simplifies the overall statistical process, from data source(s) to model application. Practically analytical techniques used in data mining include statistical methods and mathematical modeling. However, data mining and knowledge discovery is a rapidly growing area of research and application that builds on techniques and theories from many fields, including statistics, databases, pattern recognition, data visualization, data warehousing and OLAP, optimization, and high performance computing [1] . Worthy to mention that online analytical processing (OLAP) is quite different from data mining, though it provides a very good view of what is happening but cannot predict what will happen in the future or why it is happening. In fact, blind applications of algorithms are not also data mining. In particular, "data mining is a user-centric interactive process that leverages analysis technologies and computing power, or a group of techniques that find relationships that have not previously been discovered" [4] . So, data mining can be considered as a convergence of three technologies -- viz. increased computing power, improved data collection and management tools, and enhanced statistical algorithms.

Data and information have become major assets for most of the organizations. The success of any organisation depends largely on the extent to which the data acquired from business operations is utilised. In other words, the data serves as an input into a strategic decision making process, which could put the business ahead of its competitors. Also, in this era, where businesses are driven by the customers, having a customer database would enable management in any organisation to determine customer behaviour and preference in order to offer better services and to prevent losing them resulting better business. The data needed that will serve as an input to organizational decision-making process is generated and warehoused. It is being collected via many sources, such as the point of sales transactions, surveys, through the internet logs – cookies, etc. This has resulted in huge databases which have valuable knowledge hidden in them and may be difficult to extract. Data mining has been identified as the technology that offers the possibilities of discovering the hidden knowledge from these accumulated databases. Techniques such as pattern recognition and classification are the most important in data mining [4,5].

The task of recognition and classification is one of the most frequently encountered decision making problems in daily activities. A classification problem occurs when an object needs to be assigned into a predefined group or class based on a number of observed attributes, or features, related to that object. Humans constantly receive information in the form of *patterns* of interrelated facts, and have to make

decisions based on them. When confronted with a pattern recognition problem, stored knowledge and past experience can be used to assist in making the correct decision. Indeed, many problems in various domains such as financial, industrial, technological, and medical sectors, can be cast as classification problems. Examples include bankruptcy prediction, credit scoring, machine fault detection, medical diagnosis, quality control, handwritten character recognition, speech recognition etc. Pattern recognition and classification has been studied extensively in the literature. In general, the problem of pattern recognition can be posed as a two-stage process:

- **Feature selection** which involves selecting the significant features from an input pattern
- **Classification** which involves devising a procedure for discriminating the measurements taken from the selected features, and assigning the input pattern into one of the possible target classes according to some decision rule.

Research efforts dedicated to data mining, which focussed on improving the classification and prediction accuracy, have recently been undergoing a tremendous change [6,7]. The continuous development of more and more sophisticated classification models through commercial and software packages have turned out to provide some benefits only in specific problem domains where some prior background knowledge or new evidence can be exploited to further improve classification performance. In general however, related research proves that no individual data mining technique has been shown to deal well with all kinds of classification problems. Awareness of these imperfections of individual classifiers has called for the emergence of careful development and evaluation strategies of data mining classification models.

Association rule mining is a widely-used approach in data mining. Association rules are capable of revealing all interesting relationships in a potentially large database. The abundance of information captured in the set of association rules can be used not only for describing the relationships in the database, but also for discriminating between different kinds or classes of database instances.

However, a major problem in association rule mining is its complexity. Even for moderate sized databases it is intractable to find all the relationships. This is why a mining approach defines a interestingness measure to guide the search and prune the search space. Therefore, the result of an arbitrary association rule mining algorithm is not the set of all possible relationships, but the set of all interesting ones. The definition of the term interesting, however, depends on the application. The different interestingness measures and the large number of rules make it difficult to compare the output of different association rule mining algorithms. There is a lack of comparison measures for the quality of association rule mining algorithms and their interestingness measures.

Association rule mining algorithms are often compared using time complexity. That is an important issue of the mining process, but the quality of the resulting rule set is ignored. On the other hand there are approaches to investigate the discriminating power of association rules and use them according to this to solve a classification problem. This research area is called classification using association rules. It has to deal with a large number of rules.

Therefore, rule selection and rule weighting are essential for these approaches in classification. An important aspect of classification using association rules is that it can provide quality measures for the output of the underlying mining process. The properties of the resulting classifier can be the base for comparisons between different association rule mining algorithms. A certain mining algorithm is preferable when the mined rule set forms a more accurate, compact and stable classifier in an efficient way.

In the next section, we provide an overview of data mining concepts, its process, different techniques and their potential applications. In section 3, we describe our study on finding the best set of class association rules for higher predictive accuracy. Finally the paper concludes in section 4 with a glimpse to our future work.

2. Techniques and Algorithms

Researchers identify two fundamental goals of data mining: prediction and description. Prediction makes use of existing variables in the database in order to predict unknown or future values of interest, while description focuses on finding patterns describing the data the subsequent presentation for user interpretation. The relative emphasis of both prediction and description differs with respect to the underlying application and technique. There are several data mining techniques fulfilling these objectives. Some of these are classification, clustering, association and pattern discovery.

- **Classification:** Classification is the most commonly applied data mining technique, which employs a set of pre-classified examples to develop a model that can classify the population of records at large. Fraud detection and credit risk applications are particularly well suited to this type of analysis. This approach frequently employs decision tree or neural network-based classification algorithms. The data classification process involves learning and classification[2]. In Learning the training data are analyzed by classification algorithm. In classification test data are used to estimate the accuracy of the classification rules. If the accuracy is acceptable the rules can be applied to the new data tuples. For a fraud detection application, this would include complete records of both fraudulent and valid activities determined on a record-by-record basis. The classifier-training algorithm uses these pre-classified examples to determine the set of parameters required for proper discrimination. The algorithm then encodes these parameters into a model called a classifier. Some well-known classification models are:
 - a) Classification by decision tree induction
 - b) Bayesian Classification
 - c) Neural Networks
 - d) Support Vector Machines (SVM)
- **Clustering:** Clustering is a technique for identification of similar classes of objects. By using clustering techniques we can further identify dense and sparse regions in object space and can discover overall distribution pattern and correlations among data attributes. Classification approach can also be used

for effective means of distinguishing groups or classes of object but it becomes costly so clustering can be used as preprocessing approach for attribute subset selection and classification. For example, to form group of customers based on purchasing patterns, to categories genes with similar functionality. Some commonly used clustering methods are:

- a) Partitioning Methods
- b) Hierarchical Agglomerative (divisive) methods
- c) Density based methods
- d) Grid-based methods
- e) Model-based methods

2.1 Association Rules

An Association Rule is a rule of the form milk and bread \Rightarrow butter

where 'milk and bread' is called the rule body and butter the head of the rule. It associates the rule body with its head. In context of retail sales data, our example expresses the fact that people who are buying milk and bread are likely to buy butter too. This association rule makes no assertion about people who are not buying milk or bread. We now define an association rule:

Let D be a database consisting of one table over n attributes $\{a_1, a_2, \dots, a_n\}$. Let this table contain k instances.

The attributes values of each a_i are nominal. In many real world applications (such as the retail sales data) the attribute values are even binary (presence or absence of one item in a particular market basket). In the following an attribute-value-pair will be called an item. An item set is a set of distinct attribute-value-pairs. Let d be a database record. d satisfies an item set $X = \{a_1, a_2, \dots, a_n\}$ if $X \subseteq d$. An association rule is an implication $X \Rightarrow Y$ where $X, Y \subseteq \{a_1, a_2, \dots, a_n\}$, $Y \not\subseteq X$; and $X \cap Y = \emptyset$.

The support $s(X)$ of an item set X is the number of database records d which satisfy X . Therefore the support $s(X \Rightarrow Y)$ of an association rule is the number of database records that satisfy both the rule body X and the rule head Y . Note that we define the support as the number of database records satisfying $X \cap Y$, in many papers the support is defined as $s(X \cap Y) / k$. They refer to our definition of support as support count.

The confidence $\hat{c}(X \Rightarrow Y)$ of an association rule $X \Rightarrow Y$ is the fraction $\hat{c}(X \Rightarrow Y) = s(X \cap Y) / s(X)$.

From a logical point of view the body X is a conjunction of distinct attribute-value-pairs

and the head Y is a disjunction of attribute-value-pairs where $X \cap Y = \emptyset$. Coming back to the example a possible association rule with high support and high confidence would be $i_1 \wedge i_2 \Rightarrow i_3$ whereas the rule $i_1 \wedge i_3 \Rightarrow i_2$ would have a much lower support value.

2.2 Class Association Rules

The use of association rules for classification is restricted to problems where the instances can only belong to a discrete number of classes. The reason is that

association rule mining is only possible for nominal attributes. However, association rules in their general form cannot be used directly. We have to restrict their definition. The head Y of an arbitrary association rule $X \rightarrow Y$ is a disjunction of items. Every item which is not present in the rule body may occur in the head of the rule. When we want to use rules for classification, we are interested in rules that are capable of assigning a class membership. Therefore we restrict the head Y of a class association rule $X \rightarrow Y$ to one item. The attribute of this attribute-value-pair has to be the class attribute. According to this, a class association rule is of the form $X \rightarrow a_i$ where a_i is the class attribute and $X = \{a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_n\}$.

The Apriori algorithm [8,14] has become the standard approach to mine association rules. We have adapted it to mine class association rules in the way explained by Liu et al. [9,13]. The second algorithm, Predictive Apriori, has been recently proposed by Scheffer [10,12]. Both algorithms have their first step in common. They generate frequent item sets in the same way. An item set is called frequent when its support is above a predefined minimum support.

Data mining is considered to be an emerging technology that has made revolutionary change in the information world. The term 'data mining' (often called as knowledge discovery) refers to the process of analysing data from different perspectives and summarizing it into useful information by means of a number of analytical tools and techniques, which in turn may be useful to increase the performance of a system.

Technically, "data mining is the process of finding correlations or patterns among dozens of fields in large relational databases". Therefore, data mining consists of major functional elements that transform data onto data warehouse, manage data in a multidimensional database, facilitates data access to information professionals or analysts, analyse data using application tools and techniques, and meaningfully presents data to provide useful information.

According to the Gartner Group, "data mining is the process of discovering meaningful new correlation patterns and trends by sifting through large amount of data stored in repositories, using pattern recognition technologies as well as statistical and mathematical techniques"[3]. Thus use of data mining technique has to be domain specific and depends on the area of application that requires a relevant as well as high quality data.

More precisely, data mining refers to the process of analysing data in order to determine patterns and their relationships. It automates and simplifies the overall statistical process, from data source(s) to model application. Practically analytical techniques used in data mining include statistical methods and mathematical modeling. However, data mining and knowledge discovery is a rapidly growing area of research and application that builds on techniques and theories from many fields, including statistics, databases, pattern recognition, data visualization, data warehousing and OLAP, optimization, and high performance computing [1]. Worthy to mention that online analytical processing (OLAP) is quite different from data mining, though it provides a very good view of what is happening but cannot predict what will happen in the future or why it is happening. In fact, blind applications of algorithms are not also data mining. In particular, "data mining is a user-centric

interactive process that leverages analysis technologies and computing power, or a group of techniques that find relationships that have not previously been discovered" [4]. So, data mining can be considered as a convergence of three technologies -- viz. increased computing power, improved data collection and management tools, and enhanced statistical algorithms.

3. Experimental Study and Analysis

3.1 WEKA Tool

We use WEKA (www.cs.waikato.ac.nz/ml/weka/), an open source data mining tool for our experiment. WEKA is developed by the University of Waikato in New Zealand that implements data mining algorithms using the JAVA language. WEKA is a state-of-the-art tool for developing machine learning (ML) techniques and their application to real-world data mining problems. It is a collection of machine learning algorithms for data mining tasks. The algorithms are applied directly to a dataset. WEKA implements algorithms for data pre-processing, feature reduction, classification, regression, clustering, and association rules. It also includes visualization tools. The new machine learning algorithms can be used with it and existing algorithms can also be extended with this tool.

3.2 Dataset Description

We performed computer simulation on a breast-cancer dataset available UCI Machine Learning Repository [11,15]. It contains 286 samples and 9 input features as well as 1 output feature. The features describe different factor for breast-cancer reoccurrence. The output feature is the decision class which has value no reoccurrence-events and recurrence-events. The dataset contains 201 instances shown as no reoccurrence-events while 85 instances as recurrence-events. There are eight instances having missing values. A snap shot of the dataset is shown in Figure-3.1.

```

@relation breast-cancer
@attribute age {'10-19','20-29','30-39','40-49','50-59','60-69','70-79','80-89','90-99'}
@attribute menopause {'t40','ge40','premeno'}
@attribute tumor-size {'0-4','5-9','10-14','15-19','20-24','25-29','30-34','35-39','40-49'}
@attribute inv-nodes {'0-2','3-5','6-8','9-11','12-14','15-17','18-20','21-23','24-26'}
@attribute node-caps {'yes','no'}
@attribute deg-malig {'1','2','3'}
@attribute breast {'left','right'}
@attribute breast-quad {'left_up','left_low','right_up','right_low','central'}
@attribute irradiat {'yes','no'}
@attribute 'Class' {'no-recurrence-events','recurrence-events'}
@data
'40-49','premeno','15-19','0-2','yes','3','right','left_up','no','recurrence-events'
'50-59','ge40','15-19','0-2','no','1','right','central','no','no-recurrence-events'
'50-59','ge40','35-39','0-2','no','2','left','left_low','no','recurrence-events'
'40-49','premeno','35-39','0-2','yes','3','right','left_low','yes','no-recurrence-event'
'40-49','premeno','30-34','3-5','yes','2','left','right_up','no','recurrence-events'
'50-59','premeno','25-29','3-5','no','2','right','left_up','yes','no-recurrence-events'
'50-59','ge40','40-44','0-2','no','3','left','left_up','no','no-recurrence-events'
'40-49','premeno','10-14','0-2','no','2','left','left_up','no','no-recurrence-events'
'40-49','premeno','0-4','0-2','no','2','right','right_low','no','no-recurrence-events'
'40-49','ge40','40-44','15-17','yes','2','right','left_up','yes','no-recurrence-events'
'50-59','premeno','25-29','0-2','no','2','left','left_low','no','no-recurrence-events'

```

Figure-3.1: A snap shot of the dataset

3.3 Experiment Design

We follow a systematic approach to generate and prune association rules. First, we generate class association rules from the original dataset using the two selected association rule mining algorithms such as Apriori and Predicted Apriori. Using CfsSubsetEval a popular feature selection algorithm and BestFirst search, we reduced the features of the dataset in order to have better predictive accuracy. Thus the features in the two data scenarios are as follows:

Original Features: {age, menopause, tumor-size, inv-nodes, node-caps, deg-malig, breast, breast-quad, irradiat}

Reduced Features: { tumor-size, inv-nodes, node-caps, deg-malig, irradiat }

Then we apply the said association rule mining algorithms on this data scenario and derive a best set of rules for class prediction. We use confidence measure to evaluate the rules.

3.4 Results Analysis

The class association rules generated by Apriori algorithm on the original dataset is given in Figure-1 and rules generated by Predictive Apriori algorithm is shown in Figure-2.

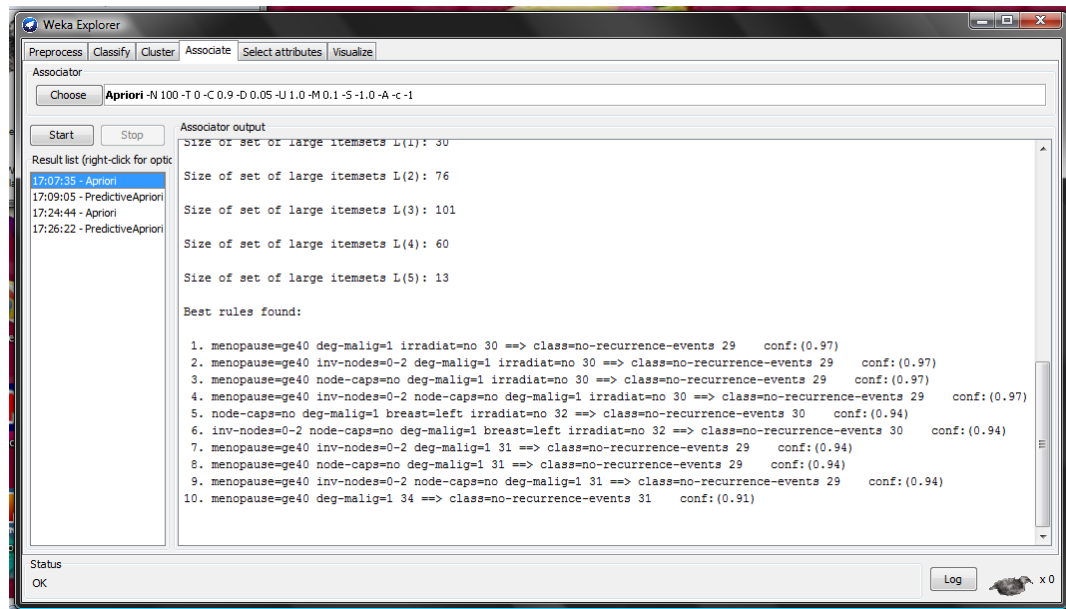


Figure-1: Rules generated by Apriori from original dataset

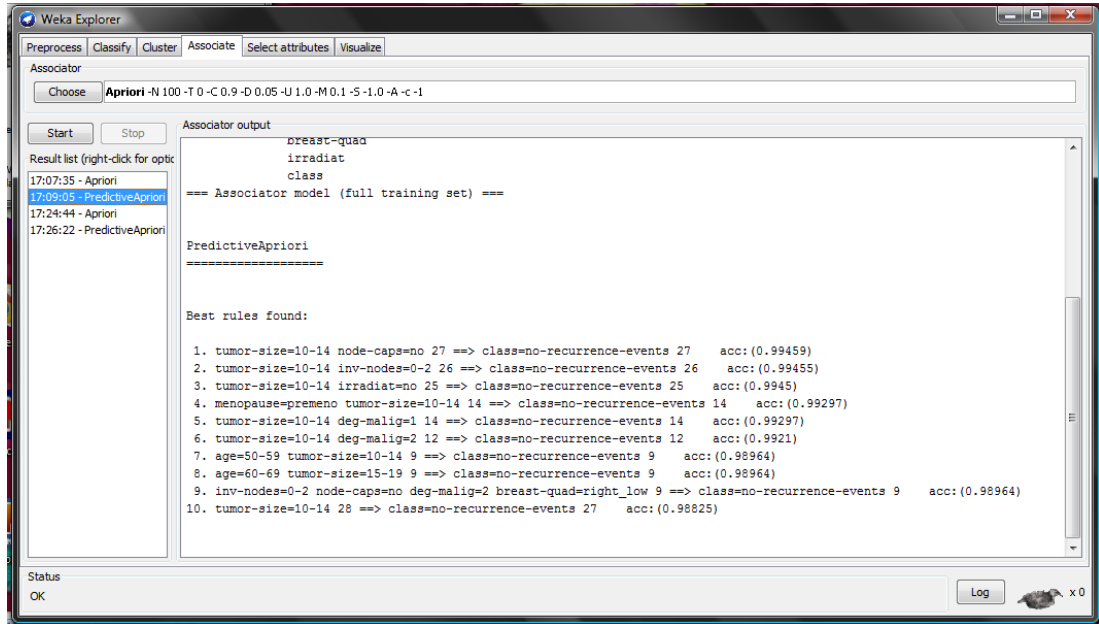


Figure-2: Rules generated by Predictive Apriori from original dataset

The class association rules generated by Predictive Apriori algorithm from the dataset with reduced features is shown in Figure-3. However, Apriori algorithm could not generate any rule form the dataset with reduced features.

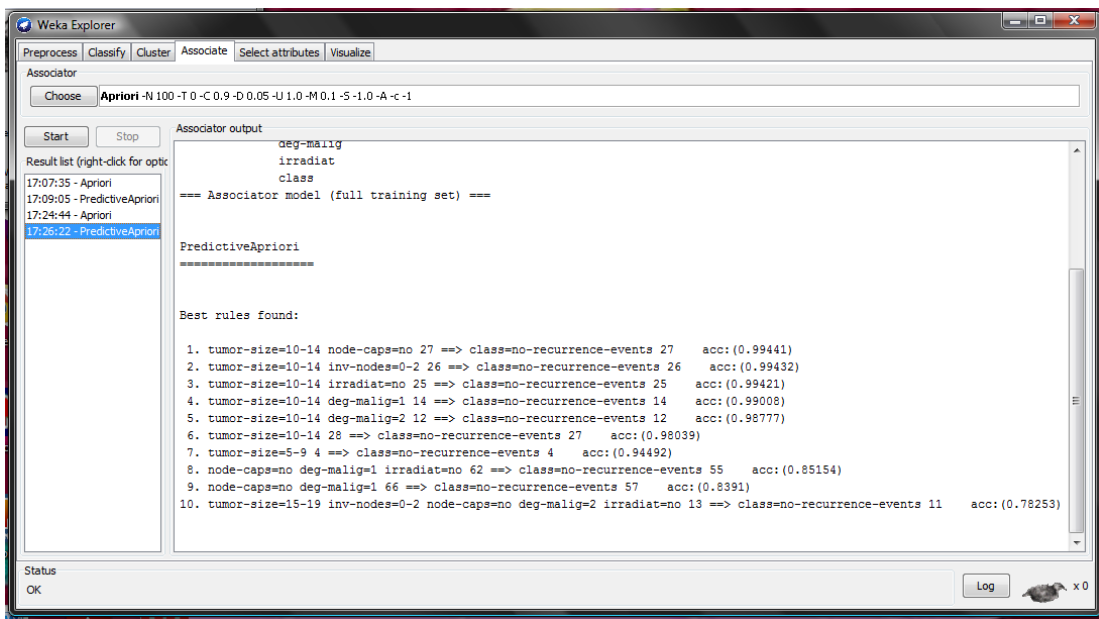
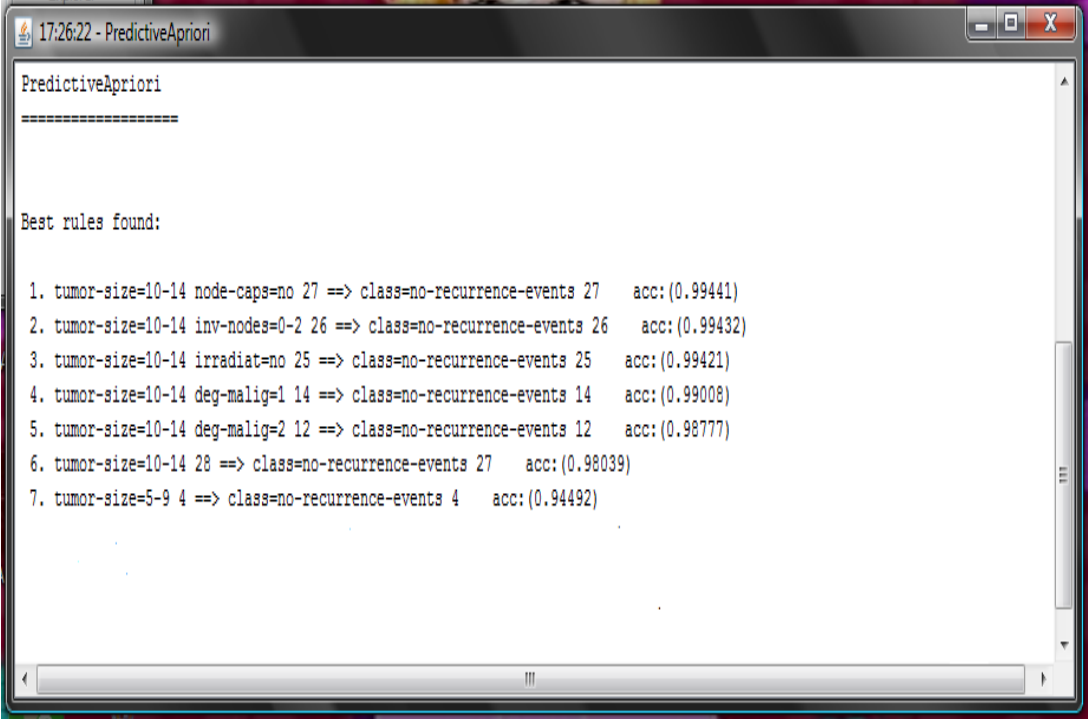


Figure-3: Rules generated by Predictive Apriori from Reduced Features

Then, we select a set of rules from the generated rule shown in Figure-3 which meet at least 90% confidence factor. These rules form the best set of rules which provide better prediction and shown in Figure-4.



```
PredictiveApriori
=====

Best rules found:

1. tumor-size=10-14 node-caps=no 27 ==> class=no-recurrence-events 27 acc:(0.99441)
2. tumor-size=10-14 inv-nodes=0-2 26 ==> class=no-recurrence-events 26 acc:(0.99432)
3. tumor-size=10-14 irradiat=no 25 ==> class=no-recurrence-events 25 acc:(0.99421)
4. tumor-size=10-14 deg-malig=1 14 ==> class=no-recurrence-events 14 acc:(0.99008)
5. tumor-size=10-14 deg-malig=2 12 ==> class=no-recurrence-events 12 acc:(0.98777)
6. tumor-size=10-14 28 ==> class=no-recurrence-events 27 acc:(0.98039)
7. tumor-size=5-9 4 ==> class=no-recurrence-events 4 acc:(0.94492)
```

Figure-4: Best set of rules for class prediction

4. Conclusion and Future Work

Data mining is a way to discover new meaning in data, performs data processing using sophisticated data search capabilities and machine learning algorithms, which can be utilized to determine the patterns or relationships implicit in a large data warehouse for better decision-making. It can be reasonably beneficial to any corporate industries, financial institutions, retailers, pharmaceutical firms, security agencies, government departments, online service providers, libraries, and individual researchers too. Business houses often use data mining to increase sales, reduce costs, improve market performance, enhance customer base by means of developing models for credit scoring, risk assessment, fraud detection, etc. In this work, we have shown how to use the classification approach to evaluate the quality of a set of association rules. In particular we applied this methodology using Apriori and Predictive Apriori to reach at a high quality set of association rules for class prediction. We conducted an experiment using WEKA data mining tool and used a dataset on health care domain from UCI machine learning repository.

4.2 Future Work

In this paper, we described a systematic approach to find the best set of association rules which perform higher predictive accuracy. We propose to extend our work by considering multiple datasets drawn from different domains to verify the findings so that the results will be sound enough for generalization.

References

- [1] Klossgen W and Zytkow J M (eds.), Handbook of data mining and knowledge discovery, OUP, Oxford, 2002.
- [2] Provost, F., & Fawcett, T., Robust Classification for Imprecise Environments. Machine Learning, Vol. 42, No.3, pp.203-231, 2001.
- [3] Larose D T, Discovering knowledge in data: an introduction to data mining, John Wiley, New York, 2005.
- [4] Kantardzic M, Data mining: concepts, models, methods, and algorithms, John Wiley, New Jersey, 2003.
- [5] Goldschmidt P S, Compliance monitoring for anomaly detection, Patent no. US 6983266 B1, issue date January 3, 2006, Available at: www.freepatentsonline.com/6983266.html
- [6] Bace R, Intrusion Detection, Macmillan Technical Publishing, 2000.
- [7] Smyth P, Breaking out of the Black-Box: research challenges in data mining, Paper presented at the Sixth Workshop on Research Issues in Data Mining and Knowledge Discovery (DMKD-2001), held on May 20 (2001), Santa Barbara, California, USA.
- [8] Agrawal R. and Srikant R. Fast Algorithms for Mining Association Rules. In M. Jarke
- [9] J. Bocca and C. Zaniolo, editors, Proceedings of the 20th International Conference on Very Large Data Bases (VLDB'94), pages 475–486, Santiago de Chile, Chile, Sept 1994 . Morgan Kaufmann.
- [10] Scheffer T. Finding Association Rules That Trade Support Optimally against Confidence. Unpublished manuscript.
- [11] Scheffer T. Finding Association Rules That Trade Support Optimally against Confidence. In L. De Raedt and A. Siebes, editors, Proceedings of the 5th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD'01), pages 424–435, Freiburg, Germany, September 2001. Springer-Verlag.
- [12] UCI Machine Learning Repository, Available at <http://archive.ics.uci.edu/ml/machine-learning-databases/statlog/german/>.
- [13] SAS Institute Inc., Lie detector software: SAS Text Miner (product announcement), Information Age Magazine, [London, UK], February 10 (2002), Available at: <http://www.sas.com/solutions/fraud/index.html>.
- [14] Berry M J A and Linoff G S, Data mining techniques: for marketing, sales, and relationship management, 2 nd edn (John Wiley; New York), 2004.

- [15] Delmater R and Hancock M, *Data mining explained: a manager's guide to customer-centric business intelligence*, (Digital Press, Boston), 2002.
- Fuchs G, *Data Mining: if only it really were about Beer and Diapers*, *Information Management Online*, July 1, (2004), Available at: <http://www.information-management.com/news/1006133-1.html>.