Level wise search algorithm based on action and non-action type data to find irregular association rules

Razan Paul¹, Abu Sayed Md. Latiful Hoque² ¹Department of Computer Science and Engineering ² Bangladesh University of Engineering and Technology Dhaka 1000, Bangladesh razanpaul@yahoo.com, asmlatifulhoque@cse.buet.ac.bd



ABSTRACT: Conventional positive association rules are frequently occuring patterns. The patterns represent what decisions are routinely made based on a set of facts. Irregular association rules are the patterns that represent what decisions are rarely made based on the same set of facts. The application domains by nature require the irregular association rules and hence we developed a level wise search algorithm that works based on action and non-action type data to find irregular association rules. Action type and non-action type will enable to discover irregular association rules. The algorithms that can tap the irregular associations have potential and our algorithm is found to be more effective in a real world datasets drawn in the health care field.

Keywords: Irregular association rules, Health care data mining, Action type data, Non-action type data

Received: 17 May 2010, Revised 27 June 2010, Accepted 4 July 2010

© 2011 DLINE. All rights reserved

1. Introduction

Data mining is applied to discover new information, which is hidden in the existing information. One of the techniques in data mining is finding association rule. The first pioneering work to mine conventional positive association rules was explained in [1]. In this work, they showed finding association rule problem can be separated into two sub problems. After that many algorithms [1-6] have been proposed to mine positive association rules efficiently. These algorithms find rules that represent decisions that occur frequently based on a set of facts. In other words, rules discovered by current association mining algorithms [1-6] are patterns that represent what decisions are usually made. In this problem, we need patterns that are rarely made. We have proposed a novel mining algorithm that can efficiently discover the association rules from the existing data that are strong candidates of variability. The algorithm uses a different candidate itemset selection process, a modified candidate generation process, and a different mechanism of generating rules from desire itemsets compared to any level-wise search algorithm. The algorithm treats all the items as being either actions that include decision, action and output or non-actions that include facts, statements and any criterion. In our problem, non-action items appear very frequently in the data, while action items rarely appear with the high frequent non-action items. Negative association rules find patterns where items are conflict with each other and do not find decision, which are made rarely. Rare association rules are the patterns containing rare items which are less frequent items and do not find decisions which are made rarely. Irregular pattern represents wrong decision, illegal practice and variability in decision.

2. Related Work

A rare association rule forms with rare data items or between frequent and rare data items. However, this rule does not map items that are used to make decision/action to antecedent and items that represent actions or decision to consequent. Moreover it does not expect antecedent to be high frequent because it find rule with high confidence. The algorithms to find

these rules only assign different support based on items frequency. A number of approaches has been proposed to mine rare associations [7-11]. In [7], minimum support is computed for each item based on the frequency of the item. In [8], A fixed minimum support is applied to extract desired itemsets, which consist of only frequent items and relative support is applied to extract desired itemsets. In [9], Negative-Binomial distribution is applied to extract Negative-Binomial desired itemsets. In [10], the association rules are found by taking into consideration only infrequent items. The approach suggested in [11] finds the association rules by computing item-wise minimum support.

A negative association rule presents a relationship among itemsets and states the presence of some itemsets in the absence of others. Every positive association rule $P \Rightarrow Q$ has three corresponding negative association rules, $P \Rightarrow \neg Q$, $\neg P \Rightarrow Q$ and $\neg P \Rightarrow \neg Q$. To extract negative association rules, most papers employ different correlation measures between attributes [12-14]. In [13], the author proposed a level-wise search algorithm for mining both positive and negative association rules that employs rule dependency measures. In [14], authors proposed another level-wise search algorithm for simultaneously extracting positive and negative association rules using Pearson correlation coefficient. In [15], author have proposed detection model using multi layer perceptron neural networks (MLP) to detect fraud/abuse problem based on medical claims. It has been proposed to detect new, unusual and known fraudulent/abusive behaviors. It works based on detection model which is very slow and need huge memory requirement to analyze existing large database. In [16], author used positive association rule to build clinical pathways, which can detect fraud and abuse on new data. However, this model cannot detect fraud and abuse from the existing large healthcare data. Our proposed approach detects fraud and abuse from the existing large information.

3. Irregular Association Rules

Let $D = \{t_1, t_2, ..., t_n\}$ be a database of n transactions with a set of items $I = \{i_1, i_2, ..., i_m\}$. Let set of action items of I be $AI = \{ai_1, ai_2, ..., ai_k\}$ where k is the number of action items. Let set of non-action items of I be $NAI = \{nai_1, nai_2, ..., nai_{m-k}\}$ where m-k is the number of non- action items. For an itemset $P \subseteq I$ and a transaction t in D, we say that t supports P if t has values for all the attributes in P; for conciseness, we also write $P \subseteq t$. By D_p we denote the transactions that contain all attributes in P. The

support of P is computed as $(P) = \frac{|D_P|}{n}$, i.e. the fraction of transactions containing P. A irregular rule is of the form: $P \models Q$, with $P \subseteq \text{NAI}, Q \subseteq \text{AI}, P \cap Q = \phi$. To hold the rule following condition must meet: P(P) or support (P) >= minimum antecedent suppot,

P(P,Q) or support (P,Q) <= maximum antecedent Consequent suppot and $\frac{P(P,Q)}{P(P)} <= maximum$ confidence where P(x) is the probability of x.



Figure 1. Data transformation of medical data

4. Mapping complex medical data to mineable items

For knowledge discovery, the medical data have to be transformed into a suitable transaction format to discover knowledge. We have addressed the problem of mapping complex medical data to items using domain dictionary and rule base as shown in figure 1. The medical data are types of categorical, continuous numerical data, Boolean, interval, percentage, fraction and ratio. Medical domain expert have the knowledge of how to map ranges of numerical data for each attribute to a series of items. For example, there are certain conventions to consider a person is young, adult, or elder with respect to age. A set of rules is created for each continuous numerical attribute using the knowledge of medical domain experts. A rule engine is used to map continuous numerical data to items using these developed rules. We have used domain dictionary approach to transform the data, for which medical domain expert knowledge is not applicable, to numerical form. As cardinality of attributes except continuous numeric data are not high in medical domain, these attribute values are mapped integer values using medical domain dictionaries. Therefore, the mapping process is divided in two phases. Phase 1: a rule base is constructed based on the knowledge of medical domain experts and dictionaries are constructed for attributes where domain expert knowledge is not applicable, Phase 2: attribute values are mapped to integer values using the dictionaries.

5. The proposed algorithm

General intuition of this algorithm is as follows: based on a set of lab tests with same results, if 99% doctors practice patients as disease x and 1 percent doctors practice patients as other diseases, then there is a strong possibility that this 1 percent doctors are doing illegal practice. In other words, if consequent C occurs infrequently with antecedent A and antecedent A occurs frequently, then $A \models C$ is a rule that is a strong candidate of variability. In every domain, there are a set of facts. Based on these facts, decision and action are taken. In a rule $S \Longrightarrow T$, if S contains a set of facts and T contains decision or action. Then such rules represent the decision T with their corresponding facts S. If $S \Longrightarrow T$ has sufficient support and confidence then it represents that decision or action T is taken routinely based on facts S. However, if S is high frequent and rule S-T has very low confidence. Then it indicates based on facts S any other decision instead of T is usually taken. It also indicates that the decision is exceptionally taken based on these facts. The main features of the proposed algorithm are as follows:

• If minimum support is only used like conventional association mining algorithm, desired itemsets that involve rarely appeared action items with the high frequent non-action items will not be found. To find rules that involve both frequent antecedent part and rare consequent items, we have used two support metrics: minimum antecedent support, maximum antecedent consequent support.

• The proposed algorithm uses maximum confidence constraint instead of widely used minimum confidence constraint to form the rules. Moreover, it partitions itemsets into action item and non-action items instead of subset generation to form rules.

• Rules have non-action items in the antecedent and action items in the consequent.

• In candidate generation, it does not check the property "Every subset of a frequent itemset is frequent" if the candidate itemset contains one or more action items to keep that itemset.

Let MAS is minimum antecedent support, MACS is maximum antecedent consequent support, I_j is the itemsets of size j, S_m is the desired itemset of size m; C_k be the sets of candidates of size k. Figure 2 shows the association mining algorithm for finding irregular rule. Like algorithm Apriori, our algorithm is also based on level wise search. Each item consists of attribute name and its value. Retrieving information of a 1-itemset, we make a new 1-itemset if this 1-itemset is not created already, otherwise update its support. The non-action 1-itemset is selected if it has support greater or equal to minimum antecedent support. The action 1-itemset is selected whatever support it has. By this way, 1-itemsets are explored which have high support for antecedent items and have arbitrary support for consequent items.

5.1 Candidate Generation

The idea behind candidate generation of all level-wise algorithms like Apriori is based on the following simple fact: Every subset of a frequent itemset is frequent so that they can reduce the number of itemsets that have to be checked. However, our proposed algorithm in candidate generation phase check this fact if the itemsets only contains non-action items. This idea

makes itemsets consist of both rare action items and high frequency non-action items. If the new candidate contains one or more action items then it is selected as a valid candidate. If the new candidate contains only non- action items then, it is selected as a valid candidate only if every subset of new candidate is frequent. This way the algorithm keeps the new candidates that have one or more action items.

5.2 Candidate Selection

We have used two separate supports metrics to filter out candidates. An itemset with only non-action items is compared with minimum antecedent support metric as non-action items can only take part in antecedent part of irregular rule, which need to be high frequent. An itemset with one or more action items is compared with maximum antecedent consequent support metric to keep rare action items with the high frequent non-action items. An itemset with only non-action items is selected if it has support greater or equal to minimum antecedent support. An itemset with one or more action items is selected if it has support greater or equal to maximum antecedent consequent support. By this way, itemsets are explored which has high support for non-action items and low support for action items with high support non-action items. Here pruning is based mostly on minimum antecedent support, maximum antecedent consequent support and checking the property "every subset of a frequent itemset is frequent".

5.3 Generating Association Rule

This problem needs association rules that represent irregular relationships between action and non-action items that occur rarely together. For this reason, the proposed algorithm uses maximum confidence constraint to form rules as it needs rule that has high support in antecedent portion and has very low support in itemset from which the rule is generated. It selects a rule if its confidence is less or equal to maximum confidence constraint. Moreover, it does not use subset generation to the itemsets to form rules. Here an itemset is partitioned into action item and non-action items. Action items are for consequent part and non-action items are for antecedent part. Here each itemset is mapped to only one rule.

5.3.1 Lemma 1. Number of rules is equal to number of desired itemsets and number of discarded rules = m^p -S where S is the number of desired itemsets. Proof: A single desired itemset consists of action type items and non-action type items. Action items and non-action items are mapped to consequent and antecedent parts respectively. Let I = { i_1, i_2, \ldots, i_n } be the set of items to be mined, where items can be either action type or non-action type. Let AI = { ai_1, ai_2, \ldots, ai_n } be the set of action items to be mined. Let NAI= { $nai_1, nai_2, \ldots, nai_v$ } be the set of non-action items to be mined. Each nai has to have confidence greater than minimum confidence support to be included as 1- itemset and all ai are included as 1- itemset. Let, C= { $c_1, c_2, c_3, \ldots, c_n$ } be the set of candidate itemsets. A new candidate NC is added to C if the non-action part of NC named NCNA holds the following property: support (each subset of NCNA) >= minimum antecedent support. A candidate c is selected for rule generation if and only if action part of c \neq NULL and c.support <= maximum confidence support. Action item and non action items of a usingle valid rule. Total rules = number of desired itemsets = S. Let m is the average number of distinct value, each multidimensional attribute holds. P is the number of attributes to be mined. Number of possible different rules = m^p . Number of discarded rules = m^p -S where S is the number of desired itemsets.

6. Results and discussion

The experiments were done using PC with core 2 duo processor with a clock rate of 1.8 GHz and 3GB of main memory. The operating system was Microsoft Vista and implementation language was c#. We used a patient dataset to verify our method. The dataset contains items, which are either actions that include decision, diagnosis and cost or non-actions that include lab tests, any symptom of patient and any criterion of disease. Each instance represents the data of one patient. We have filtered out instances which has noisy or missing values. The data set of interest has collected and preprocessed from the different local hospitals of Bangladesh, which has 50273 instances and 514 attributes (included 150 discrete and 364 numerical attributes). All these data are converted into mineable items (integer representation) using domain dictionary and rule base.

Table 1 shows test result for patient dataset, after running the program of the proposed algorithm with different parameters. Second column of the table presents the test result, where we used minimum antecedent support of 70%, maximum antecedent consequent support of 10% and maximum confidence of 10%. 49 desired itemsets were generated in total. 3 rules were discovered in total. It took about 922.2013 seconds to find these rules. Third column of the table presents the test result,



Figure 2. Association mining algorithm for finding irregular rule

where we used minimum antecedent support of 85%, maximum antecedent consequent support of 5% and maximum confidence of 5%. 31 desired itemsets were generated in total. 5 rules were discovered in total. It took about 1634.5634 seconds to find these rules.







Figure 4. Number of rules based on maximum confidence



Figure 5. Time comparison of different maximum antecedent supports



Figure 7. Accuracy of the proposed algorithm based on irregular metric



Figure 6. Time comparison of different maximum antecedent consequent supports



Figure 8. Accuracy of the proposed algorithm based on maximum Confidence

Minimum antecedent support	70%	85%
Maximum antecedent consequent support	10%	5%
Maximum confidence	10%	5%
Number of desired itemsets	49	31
Number of Desired rules	5	3
Time (Seconds)	922.2	1634.56

Table 1. Test result for patient dataset

Figure 3 shows Apriori has taken significant higher time compared to the proposed algorithm. It is because pruning in the proposed algorithm is based on minimum antecedent support, maximum antecedent consequent support and checking the property "every subset of a frequent itemset is frequent" on non-action items. Figure 4 presents if maximum confidence (MC) increases, number of valid rules increases. Figure 5 shows how time is varied with different minimum antecedent support (MAS) values for irregular rule finding algorithm. Here we measured the performance of irregular rule finding algorithm in terms of MAS keeping MACS, MC, number of action items, number of non-action items constant. Time is not varied significantly because MAS has no lead to reduce disk access as the patient data set has all sizes of candidates for these MAS values. It has only lead to the number of valid candidate generations and it can save some CPU time. As it has lead to the CPU time, the three different cases take slightly different time.

Figure 6 shows how time is varied with different MACS by keeping MAS, MC, number of action items, number of non-action items constant. Time is not varied significantly because MACS has no lead to reduce disk access as the patient data set has all sizes of candidates for these MAS values. It has only lead to the number of valid candidate generations and it can save some CPU time. As it has lead to the CPU time, the three different cases take slightly different time. As maximum consequent support decreases, number of valid candidate generation decreases For this reason, case with 5% MACS takes more time

than case with 3% MACS and case with 10% MACS takes more time than case with 5% MACS. Figure 7 illustrates accuracy results for our proposed algorithm based on minimum antecedent support. The value of minimum antecedent support for each presented result is also indicated. The figure presents MAS has no lead in accuracy, as it is not used as a parameter in selecting valid candidate and rules. Figure 8 illustrates accuracy results for our proposed algorithm based on maximum confidence. The figure presents maximum confidence has lead in accuracy as it is used as parameter in selecting valid rules. As maximum confidence decreases, accuracy increases and the number of discovered rules decreases. It is because less confidence indicates that antecedent and consequent occurs rarely together in the dataset.

7. Conclusion

Irregular patterns represent wrong decision, illegal practice and variability in decision. In this paper, we propose a level wise search algorithm that works based on action and non-action type data to find irregular association rule. The proposed algorithm has been applied to a real world patient data set. We have shown significant accuracy in the output of the proposed algorithm. Although we have used level-wise search for finding irregular patterns, each step of our algorithm is different from any other level-wise search algorithm. Rules generation from desired item sets is also different from conventional association mining algorithms.

References

[1] Agrawal, R. ImieliDski, T., and Swami, A. (1993). Mining Association Rules between Sets of Items in Very Large Databases, *In*: Proceedings of the 1993 ACM SIGMOD international conference on Management of data, Washington, D.C., p. 207-216.

[2] Agrawal, R. and Srikant, R. (1994). Fast Algorithms for Mining Association Rules in Large Databases, *In*: Proceedings of the 20th International Conference on Very Large Data Bases, San Francisco, CA, USA, p. 487-499.

[3] Brin, S., Motwani, R., Ullman, J. D., Tsur, Shalom (1997). Dynamic Itemset Counting and Implication Rules for Market Basket Data, *In*: Proceedings of the 1997 ACM SIGMOD international conference on Management of data, Tucson, Arizona, United States, p. 255-264.

[4] Mannila, H., Toivonen, H., Verkamo, A. I. (1994). Efficient Algorithms for Discovering Association Rules, *In*: AAAI Workshop on Knowledge Discovery in Databases, p. 181-192.

[5] Park, J. S., Chen, M. S., Yu, P. S. (1995). An Effctive Hash based Algorithm for mining association rules, *In*:Prof. ACM *S*IGMOD Conf Management of Data, New York, NY, USA, p. 175 - 186.

[6] Savasere, A.Omiecinski, E., Navathe, S. B. (1995). An Efficient Algorithm for Mining Association Rules in Large Databases, in Proceedings of the 21th International Conference on Very Large Data Bases, p. 432 - 444.

[7] Liu, B. Hsu, Wand Ma, Y. (1999). Mining Association Rules with Multiple Minimum Supports., *In*: SIGKDD Explorations, p. 337–341.

[8] Yun, H., Ha, D., Hwang, B., Ryu, K. H (2003). Mining association rules on significant rare data using relative support, *Journal of Systems and Software archive*, 67 (3) 181 - 191.

[9] Hahsler, M (2006). A Model-Based Frequency Constraint for Mining Associations from Transaction Data, *Data Mining and Knowledge Discovery*, 13 (2) 137 - 166.

[10] Zhou, L., Yau, S. (2007). Association rule and quantitative association rule mining among infrequent items, *In*: International Conference on Knowledge Discovery and Data Mining, San Jose, California, p. 156-167.

[11] Kiran, R. U., Reddy, P. K (2009). An improved multiple minimum support based approach to mine rare association rules, The IEEE Symposium on Computational Intelligence and Data Mining, Nashville, TN, USA, p. 340-347.

[12] Brin, S., Motwani, R., Silverstein, C (1997). Beyond Market Baskets: Generalizing Association Rules to Correlations, *In:* The Proceedings of SIGMOD, AZ, USA, p. 265-276.

[13] Wu, X., Zhang, C., Zhang, S. (2004). Efficient Mining of Both Positive and Negative Association Rules, *ACM Transactions on Information Systems*, 22 (3) 381–405.

[14] M. L. Antonie and O. R. Zaïane, "Mining positive and negative association rules: an approach for confined rules," in *Proceedings of the 8th European Conference on Principles and Practice of Knowledge Discovery in Databases*, Pisa, Italy, 2004, pp. 27 - 38.

[15] P. A. Ortega, C. J. Figueroa, and G. A. Ruz, "A Medical Claim Fraud/Abuse Detection System based on Data Mining: A Case Study in Chile," in *DMIN*, 2006, pp. 224-231.

[16] W. S. Yanga and S. Y. Hwangb, "A Process-Mining Framework for the Detection of Healthcare Fraud and Abuse.," *Expert Systems with Applications*, vol. 31, no. 1, p. 56–68, July 2006.

Conference Notification

Third International Conference on the Networked Digital Technologies (NDT 2011) University of Macau, China. July 11-13, 2011 (http://www.dirf.org/ndt)

The NDT 2011 proceedings will be published in the Communications in Computer and Information Science (CCIS) Series of Springer LNCS. The proposed conference on the above theme will be held at the University of Macau, Macau, China from July 11-13, 2011 which aims to enable researchers build connections between different digital applications. Currently, a number of institutions across the countries are working to evolve better models to provide collaborative technology services for scholarship by creating shared cyberspace thro expert collaboration, but this is a challenge for the institutions for a number of reasons. In the last few years, the landscape of digital technology applications projects for the various disciplines in humanities, social sciences, and sciences appears induced by many initiatives. For the creation of research clusters, the research community has thousands of databases, websites, local computing clusters, and web-based tools around individual themes, interests and projects. In most cases, these tools and resources are and were created to meet the specific needs of a particular community. In many cases, the funding and support for these critical initiatives is fragile and temporary, and directed in piecemeal fashion. There is a need to provide concerted efforts in building federated digital technologies that will enable the formation of network of digital technologies.

The topics to be discussed are not limited to the following:

Information and Data Management• Data and Network mining• Intelligent agent-based systems, cognitive and reactive distributed AI systems. Internet Modeling• User Interfaces, Visualization and modeling • XML-based languages• Security and Access Control• Trust models for social networks. Information Content Security• Mobile, Ad Hoc and Sensor Network Management. Web Services Architecture, Modeling and Design • New architectures for web-based social networks. Semantic Web, Ontologies (creation, merging, linking and reconciliation) • Web Services Security • Quality of Service, Scalability and Performance • Self-Organizing Networks and Networked Systems• Data management in mobile peer-to-peer networks. Data stream processing in mobile/sensor networks. Indexing and query processing for moving objects. User interfaces and usability issues form mobile applications. Mobile social networks. Peer-to-peer social networks• Sensor networks and social sensing. Social search. Social networking inspired collaborative computing. Information propagation on social networks• Resource and knowledge discovery using social networks. Measurement studies of actual social networks. Simulation models for social networks. Green Computing• Grid Computing• Cloud Computing.

All the papers will be reviewed and the accepted papers in the conference will be published in the "Communications in Com-puter and Information Science" (CCIS) of Springer Lecture Notes Series (www.springer.com/series/), and will be indexed in many global databases including ISI Proceedings and Scopus.

In addition, selected papers after complete modification and revision will be published in the following special issues journals.

- 1. Journal of Digital Information Management (JDIM)
- 2. International Journal of Information Studies (IJIS)
- 3. International Journal of Green Computing (IJGC)
- 4. International Journal of Web Applications (IJWA)
- 5. Journal of E-Technology

General Chairs

Simon Fong, University of Macau, Macau, China Sabah M.A. Mohammed, Lakehead University, Canada Sohail Asghar, Mohammad Ali Jinnah University, Islamabad

Program Chairs

Ezendu Ariwa, London Metropolitan Unversity, UK Daniel Lemire, university of Quebec at Montreal, Canada John Hamilton, James Cook University, Australia Amr Abdel-Dayem, Laurentian University, Australia Ralph Deters, University of Saskatchewan, Canada

Co-Chairs

Eric Pardede, La Trobe University, Australia Farookh Hussain, Curtin University of Technology, Australia

Workshop Chairs

Russel Pears, Auckland University of Technology, Auckland, New Zeland Sun Aixin, Nanyang Technological University, Singapore

Local Arrangement Chairs

Zhuang Yan, University of Macau, Macau, China Yang Hang, University of Macau, Macau, China

The Fourth International Conference on the Applications of Digital Information and Web Technologies (ICADIWT 2011)

Co sponsored by the IEEE Communications Society Proceedings will be published by IEEE and indexed in Xplore August 04-06, 2011 University of Wisconsin- Stevens Point USA http://www.dirf.org/diwt

This conference welcomes papers address on, but not limited to, the following research topics:

Computational Intelligence • Biometrics Technologies • Forensics, Recognition Technologies and Applications • Biometrics and Ethics • Fuzzy and neural network systems • Signal processing, pattern recognition and applications • Digital image processing • Speech processing • Computational biology and bioinformatics • Parallel and distributed computing and networks • Information retrieval and internet applications • Software engineering • Biometrics and CSR •