

ITEMSET GENERATION IN TRAINED DATA SETS USING CLASSIFICATION BASED ASSOCIATION RULE MINING

Mr. P. Pavankumar¹ & Dr. Rashmi Agarwal²

¹ Research Scholar, Computer Science and Engineering Madhav University,
Abu Road, Sirohi, Rajasthan,

² Associate Professor, Computer Science and Engineering, Madhav University,
Abu Road, Sirohi, Rajasthan.

Received: May 30, 2018

Accepted: July 04, 2018

ABSTRACT

Extracting needful information from the large pool of information is a primary task before predicting. Especially from a huge amount of incomplete, noisy, redundant and randomly scattered data. Data mining provides a framework that automatically discovers required patterns from data set which will be using these predict future occurrence in analogous scenario. Data mining approach modeled and extracts multiplicity of category and various time granularities to fulfill the need of various users or uses. Association rule mining is the core mechanism of data mining to help us. An association rule mining has grown to be central field in modern data mining research context. In this article we have generating Classification Based Association Rules (CBAR), and provide exact categorization in Classification Based Rule Mining, Classification Base Algorithm provides exact trained datasets various techniques of association rule mining and their significance.

Keywords: Association Rule Mining (ARM), CBAR, Classification Base Algorithm

I. INTRODUCTION

Extracting needful information from the large pool of information is a primary task before predicting. Especially from a huge amount of incomplete, noisy, redundant and randomly scattered data. Here Data mining plays a significant role to specify the model which extracts the knowledge from the large information which will be the decisive next (data). Broadly Data mining approach can be classify into two description and prediction. In Descriptive type mining characterization of common distinctiveness of data while Predictive is to predict on the basis of the referring of data.

Data mining provides a framework that automatically discovers required patterns from data set which will be using these predict future occurrence in analogous scenario. Data mining approach modeled and extracts multiplicity of category and various time granularities to fulfill the need of various users or uses. The author of [31] divides data mining tasks in following (a) conceptual description; (b) correlation analysis; (c) classification and prediction; (d) clustering; (e) outlier analysis; (f) evolution analysis. Development and Éclat calculations. A different methodology either improves the productivity of the current methodologies or manages abnormal state data mining ideas. The most agreeable order of data mining techniques is on the premise of the design of the database under thought. Diverse methodologies have been suggested that utilization even design of

database, vertical format of database or anticipated format of database. A few scientists deal with enhancing the productivity of the mining process while others attempted to uncover progressed, confused and abnormal state information from the database. Additionally, swarm insight techniques have been utilized as a part of various fields for different assignments going from advancement to appropriation of assets. The utilization of swarm insight for data mining has turned out to be well known since most recent two decades. After that few developments in the field of data mining utilizing swarm knowledge has been completed. This section contemplates light on the accessible writing in both the field's viz. data mining and swarm knowledge and likewise introduces a talk of the fruitful applications of various swarm insight techniques in data mining.

1.1 Association Rule Mining

Data and Author of [9] ha illustrated the working of Association rule mining in article, According to the author defined as discovers the interesting associations among items in a given data set. Suppose a Database T is a collection of m transactions, $\{T_1, T_2, \dots, T_m\}$ and I is the set of all items, $\{i_1, i_2, \dots, i_n\}$, where each of the transactions T_k ($1 \leq k \leq m$) in the database T represents a set of items ($T_k \subseteq I$). An item set is defined as a non-empty subset of I.

Now the association rule defined as: $X \rightarrow Y(c, s)$, where $X \subseteq I$, $Y \subseteq I$ and $X \cap Y = \emptyset$.

Where s = support and
 c = confidence

The support calculated as is the percentage of the transactions in which both variable X and Y is appeared in the similar transaction while confidence is the proportion of the number of transactions which contains both X and Y to the number of transactions that contain only X .

It can be formulated as follows:

Support $(X \rightarrow Y) = P(X \cup Y)$

Confidence $(X \rightarrow Y) = P(Y|X)$

Association rule mining is to help find the relationship between the item sets in a large number of databases. By describing the potential rules between the items in databases, dependencies between multiple domains which meet the given support and the confidence threshold are found. The relationship is hidden and unknown in advance, is not got by the logic operation of a database or statistical methods. They are not the inherent properties of the data itself, but on the characteristics of data items simultaneously. The most typical example of association rules cases is: "80 percent of customers buy beers also buy diapers at the same time", its intuitive meaning is, how larger the tendency of customers buy certain products while they will buy other goods. Discovering the rules like this is valuable for setting marketing strategy. Association rules can be applied to analysis of customer shopping, product shelf design, mailing of commercial advertising, catalog design, additional sales, storage planning, network fault analysis, and classification of the users based on buying patterns.

Association rules mining is to find out all the strong association rules in transaction database D with user-given minimum support and minimum confidence. Corresponding itemsets of strong association rules $A \rightarrow B$ must be frequent itemsets, and the confidence of association rules $A \rightarrow B$ derived from frequent itemsets A B is calculated by the frequent itemsets A and $A \cup B$'s support. Therefore, the association rules mining can be decomposed into two steps: The first step is to find out all the frequent itemsets in D quickly and efficiently, which is the central issue of association rules mining.

The second step is to produce a strong frequent itemsets. We use frequent itemsets to generate the required association rules, based on the user-given minimum confidence to select. Finally, strong association rules are generated.

Agrawal et al. has attest the association rules have the following characteristics-

1. The subset of the frequent items is also frequent items.
2. The superset of the non-frequent items is also non-frequent.

II. CLASSIFICATION BASED ASSOCIATION RULES (CBA)

The rules coming about because of Associative Classification mining can be assessed to choose a subset of the rules that will frame the model or classifier. To the best of our insight, Liu, Hsu, and Ma were the first to create a classifier in light of affiliation rules. They demonstrate that the classifier constructed executes and also superior to anything surely understood decision tree calculations. From that point forward, numerous affiliation administer based classifiers have been worked for different areas. Among others, for classifying mammography pictures, for classifying web documents, for recommender frameworks, for classifying spatial information, for document classification, and for content arrangement. The way toward building the classifier includes choosing rules by certainty or support. Certainty is a well known standard for run determination to the classifier as it signifies the strength of a run the show. On account of CBA, they utilize a heuristic to choose a subset of the rules that orders the preparation set generally precisely. At times, the pruning is as basic as evacuating contradicting rules or more confused like utilizing post pruning techniques that are utilized as a part of decision trees.

2.1 Rule Ordering Association

Given two rules, r_i and r_j , r_i goes before r_j if the certainty of r_i is more noteworthy than that of r_j or, their certainty are the same, however the help of r_i is more prominent than that of r_j , or, both the certainty and the help of r_i and r_j are the same, yet r_i is produced sooner than r_j .

Give R a chance to be the arrangement of CARs and D be the preparation information. The point of the model development algorithm is to pick an arrangement of profoundly predictive rules in R to cover the preparation information D . The classifier constructed is of the accompanying structure: $\langle r_1; r_2; \dots; r_n; \text{default class} \rangle$ where $r_i \in R$, $r_a r_b$ if $a < b$. Default class is the default mark utilized when none of the rules can classify a case. Algorithm 2 demonstrates the CBA-CB technique. In stage 1, the rules are sorted by the request said above; at that point each administer is considered thusly.

Algorithm: CBA-CB Algorithm

Inputs: rules R , training set instances D

Output: classifier C

14. $R = \text{sort}(R)$;

15. for each rule $r \in R$ in sequence do

16. temp = ;

17. for each instance $d \in D$ do
18. if d satisfies the conditions of r then
19. store $d.id$ in $temp$ and mark r if it correctly classifies d ;
20. end if
21. end for
22. if r is marked then
23. insert r at the end of C ;
24. delete all the cases with the ids in $temp$ from D ;
25. select the default class for the current C ;
26. compute the total number of errors of C ;
27. end if
28. end for
29. Find the rule p in C such that C_p , the list of rules in C up to p , has the lowest total number of errors. and drop all the rules.
30. Add the default class associated with p to the end of C , and return C

The rule under thought is checked on the off chance that it can classify no less than one occurrence in the training set effectively (stages 5 and 6). In the event that the rule is denoted, every one of the instances secured by the rule are expelled from the training set and the lion's share class of whatever remains of the training instances turns into the default class mark (stages 11 and 12). The stamped rule is added to the finish of the classifier.

Give C_r a chance to indicate the arrangements of rules finishing off with rule r that have been chosen for inclusion in the classifier up until now. In stage 13, the classifier C_r is utilized to classify the instances of the training set, and assess the execution of the classifier. Since the classification estimations of the instances are known, every classification endeavor or forecast can be recorded as a right classification or wrong classification. At the point when every one of the instances are ordered, the classifier will be appointed an error rate which is the aggregate number of wrong classifications over the aggregate number of classifications.

The rule for which C_p has the least number of errors is found and all rules included after this rule is evacuated. The default class name connected with that rule turns into the default class name of the classifier

2.2 PostClassificationAssociation Rules

With association rule mining, the quantity of rules created may overpower. As all the created rules may not intrigue or huge, it is essential to prune those rules regarded uninteresting or over fitting (rules that are extremely specific to the training set). Like decision tree post-pruning,

associations rules can be present pruned on diminish the quantity of rules delivered. Numerous thoughts on post-pruning of decision trees were represented by Quinlan. There are essentially two approaches to post-pruning in view of error rates. One is to separate the informational index into training, validation and testing sets. With this approach, the rules will be constructed utilizing the training set, and pruning will be done in light of the execution of the rules on the validation set. With the second approach, there is no different validation set, yet the training set is utilized as the validation set. The last technique is known as pessimistic error pruning.

Pessimistic error pruning is a heuristic in view of statistical thinking. For each rule, let the quantity of errors on the training set be E and the quantity of cases secured on the training set be N (those instances containing the predecessor of the rule). The watched error rate is $f = E/N$. Give the genuine error a chance to rate (obscure) be q . Here, we accept the N instances are produced by a Bernoulli procedure with probability q and error rate E .

The mean and variance of a solitary Bernoulli trial with progress rate p will be p and $p(1-p)$ individually. For N Bernoulli trials, the achievement rate f is an arbitrary variable with mean equivalent to p , and the variance is diminished to $p(1-p)/N$. For expansive N , the estimation of the arbitrary variable f approaches an ordinary distribution.

The probability that a random variable, X , with 0 mean lies within a con dense range of width $2z$ is

$$Pr [-z \leq X \leq +z] = c$$

Where c is the confidence level.

For random variable f to have a 0 mean and unit variance, we subtract mean p from f and divide by standard deviation,

where

$$\sigma = \sqrt{p(1-p)/N}$$

$$Pr \left[\frac{f - p}{\sqrt{p(1-p)/N}} > z \right] = c$$

The upper confidence limit for q in the expression above provides a pessimistic estimate of e error rate at a given node:

$$e = \frac{f + \frac{z^2}{2N} + z \sqrt{\frac{f}{N} - \frac{f^2}{N} + \frac{z^2}{4N^2}}}{1 + \frac{z^2}{N}}$$

Rule R is compared with its subrules, which are rules in which one or more items are removed from the antecedent of R. If a rule R has a higher pessimistic error rate than any of its subrules, R is pruned while retaining the subrules.

III. Association Rule Based Classification Model

Data have been acknowledged as a profitable resource since long time. In any case, the use of data and the apparatuses for utilizing that data has been changed a considerable measure after some time. During 1960's data Data and data have been acknowledged as a profitable resource since long time. In any case, the use of data and the apparatuses for utilizing that data has been changed a considerable measure after some time.

3.1 Classification Association Rules Generation

Each CAR has a class property or focus on the ensuing of the rule. As this objective is predefined, we can utilize this objective as a semantic limitation to create visit thing sets comprising of the class property.

Definition 3.1.1. Semantic Constraints

A semantic constraint is a requirement that an attribute(s) must appear or must not appear in the antecedent and/or consequent of a rule.

Definition 3.1.2. Syntactic Constraints

A Syntactic constraint is a requirement placed on the number of attribute-value pairs on either the antecedent or consequent of a rule.

3.2 Generating Rules with Semantic Constraints

These three attributes add to the semantic constraints. For the rules to have these three attributes, the incessant item sets must contain them. In this way, it will source in the event that we create just things sets that incorporate the three attributes we are occupied with. The regular item set may have different attributes. We can utilize the semantic constraints as conditions in the join advance of the Apriori competitor generation stage.

The approach we have used to prune item sets that don't contain the required attributes is firmly identified with the implementation of Apriori Sets And Sequences. We will probably create just item sets that have all the required attributes (constraints).

In the Apriori Sets and Sequences algorithm, each quality esteem combine is mapped to a number (thing number). This numbering is done such that the characteristic estimations of the rst trait in the informational collection gets the most minimal numbers took after by the property estimations of the second quality et cetera. A hash table stores the mapping between the numbers and the quality

esteems. Numbers appointed to a property's estimations are consecutive.

To consider pruning of item sets that may not contain the attributes we want, we reorder the attributes with the goal that the attributes that are semantic constraints (required attributes) are given littler numbers than the non-required attributes. In this way, in the contact-focal point's informational index, characteristic estimations of contacts, age, and tear goad rate will be allotted littler numbers than the estimations of the other quality, astigmatism, as show beneath.

3.2.1 Generating Frequent Item sets

Table 3.1: Attribute values renumbered to give lower numbers

contact-lenses=soft	1
contact-lenses=none	2
contact-lenses=hard	3
age=young	4
age=pre-presbyopic	5
age=presbyopic	6
tear-prod-rate=normal	7
tear=prod-rate=reduced	8
astigmatism=yes	9
astigmatism=no	10

Let us generate frequent itemsets from the contact-lenses data set:

In Figure3.1, we showed the attribute-value pairs and how they are numbered.

The required attributes are: contacts, age and tear prod rate.

Association rule mining is one of the essential and all around investigated techniques of data mining to discover vital connections among data things. In light of the design of database, 28 distinct techniques have been created for mining the data. The flat design of database is utilized by Apriori arrangement strategies while vertical format is the base of FP-development and Éclat calculations. A different methodology either improves the productivity of the current methodologies or manages abnormal state data mining ideas. These all around preferred methodologies are talked about below:

In Tables 3.1, 3.2 and 3.3 we show the candidate itemsets generated and their support until no more candidate itemsets can be generated. We use minimum support as 1, at

least one data instance must contain the itemsets for the itemsets to be considered for the next level. Those itemsets with support less than 1 were

Table3.2: Candidate itemsets in the second level of itemset generation

itemset	Support	itemset	Support	itemset	Support	itemset	Support
f1, 2 g	0	f1, 3g	0	f1, 4g	1	f1, 5g	1
f1, 6g	0	f1, 7g	2	f1, 8g	0	f1, 9g	0
f1, 10g	2	f2, 3g	0	f2, 4g	1	f2, 5g	2
f2, 6g	3	f2, 7g	2	f2, 8g	4	f2, 9g	3
f2, 10g	3	f3, 4g	1	f3, 5g	1	f3, 6g	1
f3, 7g	3						

Table 3.3: Candidates item sets in the third level of itemset generation

itemset	Support	Itemset	Support	itemset	Support
f1, 4, 5g	0	f1, 4, 7g	1	f1, 4, 10g	1
f1, 5, 7g	1	f1, 5, 10g	1	f2, 4, 5g	0
f2, 4, 6g	0	f2, 4, 7g	0	f2, 4, 8g	0
f2, 4, 9g	0	f2, 4, 10g	0	f2, 5, 6g	0
f2, 5, 7g	1	f2, 5, 8g	1	f2, 5, 9g	1
f2, 5, 10g	1	f3, 4, 5g	0	f3, 4, 6g	0
f3, 4, 7g	1	f3, 4, 9g	1	f3, 5, 6g	0
f3, 5, 7g	1	f3, 5, 9g	1	f3, 6, 7g	1
f3, 6, 9g	1				

3.2.2 Generating Maximal Frequent Item sets

All the frequent item sets produced from the item set generation step are utilized to create maximal frequent item sets. A maximal frequent item set is an item set that isn't a subset of some other item set. This stage reduces the quantity of item sets we are working with essentially.

3.2.3 Counting Support for those Item sets without Support

Utilizing the maximal item sets, we produce every one of the subsets of the maximal item sets and decide whether every subset has support tallied. As we probably am aware from the item-set generation organize, some item sets may not be produced in light of the item set pruning step and thusly won't have their support tallied. Those item sets without support should have their

support checked before the rule generation organizes. While generating rules from the maximal item sets, a pruned item set may show up on the precursor or resulting of a rule. Assume X speaks to the predecessor and Y speaks to the subsequent, confidence of a rule is registered as supportfX [Y g supportfXg]. In this manner, it is essential that all subsets of a maximum minimal item set have their support checked.

IV CONCLUSION

Our results demonstrate that utilizing item set pruning within the sight of constraints dewrinkles the quantity of resulting maximal item sets and additionally enhances the time taken for the procedure of item set mining. In any case, in our present implementation the general mining time (i.e., the item set mining time in addition to

the rule development time) may take longer than the general mining time when non-item set pruning is used.

Additionally work should be done around there to enhance the rule development process so the rule generation time does not overshadow the time reserve funds obtained by item set pruning during item set generation.

REFERENCES

1. Xiufeng Piao, Zhanlong Wang and Gang Liu –Research on Mining Positive and Negative Association Rules Based on Dual Confidence||, IEEE, Research on Mining Positive and Negative Association Rules Based on Dual Confidence, 2010.
2. Agrawal R, Imielinski T, Swami A, –Mining Association Rules between Sets of Items in Large Databases, ||In: Proc of the ACM SIGMOD International conference on Management of Data, Washington DC, 1993, pp, 207-216.
3. Agrawal R, Srikant R, –Fast algorithms for mining association rule, ||In: Proc. 20th Int. Conf. on VLDB. Santiago, Chile, 1994, pp, 487-499.
4. Mojdeh Jalali-Heravi, Osmar R. Zaiane, –A Study on Interestingness Measures for Associative Classifiers, || TProceedings of the 2010 ACM Symposium on Applied Computing, 2010, pp, T1039-1046.
5. He Jiang, Yuanyuan Zhao, Chunhua Yang, Xiangjun Dong, –Mining Both Positive and Negative Weighted Association Rules with Multiple Minimum Supports,|| 2008 International Conference on Computer Science and Software Engineering, 2008, pp, 407-410.
6. Ling Zhou, Stephen Yau, –Association Rule and Quantitative Association Rule Mining among Infrequent Items,|| Proceedings of the 8th international workshop on Multimedia data mining: associated with the ACM SIGKDD 2007.
7. H Liqiang, GengH, HHoward J. HamiltonH, –Interestingness measures for data mining: A survey.|| ACM Computing Surveys , 2006, pp, 1-32.
8. SHANG Shi-ju, DONG Xiang-jun, LI Jie, –Algorithms for mining negative association rules in multi-database, ||Computer Engineering and Applications, 2009, 45(24), pp, 150-152.
9. [KaiXing Wu, Juan Hao and Chunhua Wang –Application of Fuzzy Association Rules in Intrusion Detection||, IEEE, International Conference on Internet Computing and Information Services, 2011.
10. Ruowu Zhong and Huiping Wang –Research of Commonly Used Association Rules Mining Algorithm in Data Mining||, IEEE, International Conference on Internet Computing and Information Services, 2011.
11. Agrawal R, Imielinski T, Swami A. Mining Association Rules between Sets of Items in Large Databases [C] // Proceedings of the 1993 ACM SIGMOD Conference. Washington D C, 1993: 207-216
12. Savasere A, Omieeinski E, Navathe S. An efficient algorithm for mining association rules in large databases[C]//Proceedings of the 21st International Conference on Very Large Databases. Zurich, Switzerland, 1995: 432-443
13. Brin S, Motwani R, Ullman J D. Dynamic Itemset counting and implication rules for market basket data[C]//Proceedings of the 1997 ACM SIGMOD International Conference on Management of Data. Tucson, Arizona, 1997: 255-264
14. Hart J, Pei J, Yin Y. Mining frequent patterns without candidate generation[C] //In Proceeding of the 2000 ACM SIGMOD Conference on Management of Data. Dallas, TX, 2000 [6] Liu J, Pan Y, Wang K, et al. Mining frequent item sets by opportunistic projection[C]//Proc Of the Eighth ACM SIGKDD Intl. Conf on Knowledge Discovery and Data Ming. Alberta, Canada, 2002: 229—238
15. Agarwal R, Aggarwal C, Prasad V V V. A tree projection algorithm for generation of frequent itemsets [J]. Parallel and Distributed Computing, 2001, 61(3): 350-371
16. Agrawal R, Srikant R. Fast Algorithm for Mining Association Rules [C] //Proceedings of the 20th Very Large Data Bases (VLDB '94) Conference. Santiago, Chile, 1994: 487-499
17. Agrawal R, Srikant R. Fast Algorithm for Mining Association Rules in Large Databases[R]. San Jose, IBM Alma den Research Center, 1994