

Dynamic System to Discover a Pattern

S.Thabasu Kannan¹ and J.Malarvizhi²

¹Principal, Pannai College of Engineering & Technology, Sivagangai – 630 561, Tamil Nadu, India

²Research Scholar, Bharathiar University, Coimbatore. Tamil Nadu, India

E-Mail: thabasukannan@gmail.com, malar76@gmail.com

Abstract - It is indeed an art to match maximum number of preferences by utilizing limited number of resources. During the current academic year 75% of the admissions to Engineering Colleges have gone down, as only 30% to 40% of intake has been filled. Without reaching the breakeven point, the management of the institution becomes a complicated issue. In this situation providing quality education to the students is the question mark. The main aim of this paper is to discover a pattern to identify the choice of preferences of the candidates to seek admissions in any academic institutions. The candidate finds admission in an institution only when his/her own preference matches exactly, otherwise the candidate continues to go by the next alternate in the list of preference. If the institution analyzes the preferences of the candidates and tries to satisfy them, surely the institution can reach even above their intake. Generally satisfaction of individual candidate is practically not possible. Hence the institution should try to satisfy maximum number of candidates by utilizing our existing infrastructure and viable number of preferences. Here the viability is the main constraint. For the purpose of matching optimum number of candidates to suit our existing system, we have designed our algorithmic approach. Here our new system is used to extract frequent item sets from various preferences. By thresholds, it can fix the preferences either decrease or increase the level of frequent. The new algorithm is based on association rule classification which is one of data mining techniques. Data mining is the process of extracting knowledge hidden from large volumes of raw data.

It is based on the concept of prune. Here the frequency of itemset₂ is combined with frequency to get itemset₃ and continues until itemset_n. The new algorithm is easy to use and implement because its complexity is less. The application is designed to generate association rule until n-antecedent with one consequent. For this study purpose we have identified 15 most frequently used preferences among the students. The samples we have taken to get association rules are 100 students of Pannai College of Engineering and Technology at Sivagangai. The discovered pattern is common to all institutions. The pattern discovery may be accurate because it is computed by using factors like confidence and support. If this intelligent system is followed strictly, definitely the number of outcomes is increased. The applicant would prefer only when the supply is high. The result of this paper is an application that can generalize association rule among various academic institutions.

I. INTRODUCTION

In order to avoid the problem of processing big data into information which useful for user, we can use data mining techniques. Data mining [6-9] is the analysis of data with the intent to discover gems of hidden information in the vast

quantity of data that has been captured in the normal course of running the business.

Here we develop an application to classify students based on their preferences of academic institutions in one transaction using our new algorithm. We show how to extract data pattern with assumption that data has been in one warehoused database.

The expected outcomes of this paper will be the patterns of student's preferences to study that are bought in one transaction together. These patterns can be used to get admissions for any academic institutions that frequently received in a same time, to design admission diary, to design a template, to plan to get admissions etc.

II. NEW ALGORITHM

```
{
L1 = {large l-itemsets};
for ( k = 2; Lk-1 ≠ 0; k++ ) do begin
ck = thabu-gen(Lk-1);
for all transactions t ∈ D, do begin
Ct = subset,(Ck, t);
For all candidates c ∈ Ct, do
c.count++;
end
Lk; = {c ∈ Ck | c.count ≥ MinSup}
end
Answer = UkLk;
}
```

The newly developed algorithm is for mining association rules, takes advantage of structure within the rules themselves to reduce the search problem to a more manageable size. Suppose that an academic institution is having the following facilities {moderate fees, state of the art lab, high volumes of library, experienced faculty, good result, placement, hostel facility, transport facility, sports activity, basic amenities, hi-tech environment, good management, reputation, distance, reference}.

Denote this set of items as I. One by one, the candidates pull over, pick up their preferences and get admissions through various combinations of these items, subsets of I. Suppose the following table, lists the transactions made during the academic year 2014-15.

TABLE 1 TRANSACTIONS MADE DURING THE ACADEMIC YEAR 2014-15

S. No.	Primary Item (A)	Secondary Items (B)	No. of Trans. A	No. of Trans. A & B	Support	Confidence
1	Moderate fees	Placement, Good result	8	5	0.074	0.625
2	State of the art lab	Moderate fees , placement	8	6	0.089	0.75
3	High volume Library	Moderate fees, good result	4	3	0.044	0.75
4	Experienced faculty	Placement, State of the art lab	5	4	0.060	0.8
5	Good result	State of the art lab, Placement	9	5	0.074	0.56
6.	Placement	Moderate fees, Experienced faculty	8	5	0.074	0.625
7	Hostel facility	Moderate fees, Transport facility	3	1	0.015	0.33
8	Sports activity	Experienced faculty, Placement	3	2	0.03	0.67
9	Basic amenities	State of the art lab, Transport facility	2	1	0.015	0.5
10	Transport facility	Placement, Basic amenities	4	2	0.03	0.5
11	Hi-tech environment	Placement, Transport facility	3	2	0.03	0.67
12	Good management	Placement, Library	2	1	0.015	0.5
13	Reputation	Experienced faculty, State of the art lab	3	2	0.03	0.67
14	Distance	Experienced faculty, Transport facility	2	1	0.015	0.5
15	Reference	Moderate fees, State of the art lab	3	1	0.015	0.33
Total			67	44		

Let D be the set of transactions represented in the Table1, where each transaction T in D represents a set of items contained in I . Suppose that we have a particular set of items A (e.g., moderate fees and Placement), and another set of items B (e.g., moderate fees and state of the art lab). Then an *association rule* takes the form *if A, then B* (i.e., $A \rightarrow B$), where the *antecedent A* and the *consequent B* are proper subsets of I , and A and B are mutually exclusive. This definition would exclude trivial rules such as *if moderate fees and state of the art lab then Placement*.

$$Support = \frac{No\ of\ transactions\ containing\ both\ A\ \&\ B}{Total\ no\ of\ transactions}$$

$$confidence = \frac{no\ of\ transactions\ containing\ both\ A\ \&\ B}{no\ of\ transactions\ containing\ A}$$

The mining of association rules from large databases is a two-steps process (74,78):

- Find all frequent itemsets i.e find all itemsets with frequency $\geq \Phi$.
- From the frequent itemsets, generate association rules satisfying the minimum support and confidence conditions.

a. Generating Frequent Itemsets

This step is looking for items combination with frequency $\geq \Phi$. The global candidate set is generated as the union of all local large itemsets from all partitions. It also takes n (no of partitions) iterations. The *support* for a candidate itemset in that partition is generated by intersecting all l-subsets of that itemset. The cumulative count gives the global support for the itemsets.

An *itemset* is a set of items contained in I , and a *k-itemset* is an itemset containing k items. For example, {moderate fees and placement} is a 2-itemset, and {experienced faculty, good result, State of the art lab} is a 3-itemset, each from the set I . The *itemset frequency* is simply the number of transactions that contain the particular itemset. A *frequent itemset* is an itemset that occurs at least a certain minimum number of times. We denote the set of frequent k -itemsets as F_k . We first find F_1 , the frequent 1-itemsets, which represent simply the individual items themselves. Since each sum meets or exceeds $\Phi = 3$ (threshold), we conclude that each 1itemset is frequent.

In general, to find F_k , our new algorithm first constructs a set C_k of candidate k -itemsets by joining F_{k-1} with itself. Then it prunes C_k using the a priori property. The itemsets in

C_k that survive the pruning step then form F_k . Here, C_2 consists of all the combinations of items. Here no of transactions A&B is called candidates 2-itemset (F_2).

Since $\Phi = 3$, we have $F_2 = \{\{Placement, good result\}, \{Moderate fees, Placement\}, \{Moderate fees, Good result\}, \{Placement, State of the art lab\}, \{Moderate fees, Experienced faculty\}\}$

Next, we use the frequent itemsets in F_2 to generate C_3 , the candidate 3-itemsets. To do so, we join F_2 with itself, where *itemsets are joined if they have the first $k-1$ items in common*. For example, $\{Moderate fees, Placement\}$ and $\{Moderate fees, Good result\}$ have the first $k-1 = 1$ item in common, Moderate fees. Thus, they are joined into the new candidate itemset $\{Moderate fees, Placement, Good result\}$. Similarly, $\{Placement, good result\}$ and $\{Placement, state of the art lab\}$ have the second item, Placement in common, generating the candidate 3-itemset $\{Placement, good result, state of the art lab\}$. Finally, candidate 3-itemsets $\{moderate fees, placement\}$, $\{Placement, State of the art lab\}$ are generated in like fashion. Thus, $C_3 = \{\{Moderate fees, Placement, State of the art lab\}, \{Moderate fees, Good result, Experienced faculty\}\}$.

C_3 is then pruned, using our newly developed algorithm. For each itemset s in C_3 , its size $k-1$ subsets are generated and examined. If any of these subsets are not frequent, s cannot be frequent and is therefore pruned. For example, let $s = \{Moderate fees, Placement, State of the art lab\}$. The subsets of size $k-1 = 2$ are generated, as follows: $\{Moderate fees, Placement\}$, $\{Moderate fees, State of the art lab\}$, and $\{Placement, State of the art lab\}$.

From the above, we see that each of these subsets are frequent and that therefore $s = \{Experienced faculty, Placement, Well library\}$ is not pruned. Because here the subset $\{Experienced Faculty, Transport facility\}$ has frequency $1 < 3 = \Phi$ so that it is not frequent. By our new algorithm, $\{Experienced Faculty, Placement, Transport facility\}$ cannot be frequent, is therefore pruned, and does not appear in F_3 . Also consider $s = \{Moderate fees, Transport facility, State of the art lab\}$. The subset $\{Moderate Fees, Transport facility\}$ has frequency $1 < 3 = \Phi$, and hence is not frequent. Again, by the defined property, its superset cannot be frequent and is also pruned, not appearing in F_3 .

b. Generating Association Rule

After all of the frequent itemset has been found, the next step is generating association rule [74, 79] by using confidence formula. Once the large itemsets and their supports are determined the rules can be discovered in a straight forward manner as follows: if I is a large itemset, then for every subset a of I, the ratio support (I) / support (a) is computed. If the ratio is at least equal to the user specified minimum confidence, then the rule $a \Rightarrow (I - a)$ is output.

The local large itemsets [75, 80] are generated for the same minimum support as specified by the user. Hence this is equivalent to generate large itemsets with that minimum support. For large partition sizes, the number of local large itemsets is likely to be comparable to the number of large itemsets generated for the entire database.

From $F_3: \{Moderate fees, Placement, Good result\}$, we have 3 candidate rules that have 2 antecedent and one consequent.

TABLE 2 CANDIDATE RULES OF F3

S.No	Rule	Confidence	
1	if <i>Moderate fees</i> and <i>Placement</i> Then <i>Good result</i>	6/9	67%
2	if <i>Moderate fees</i> and <i>Good result</i> Then <i>Placement</i>	3/8	37.5%
3	if <i>Placement</i> and <i>Good result</i> then <i>moderate fees</i>	5/8	63%

If our minimum confidence is 65% then the second and third rule would not be reported. Finally we turn to single antecedent and single consequent. The candidate rule is shown in the below table:

TABLE 3 CANDIDATE RULES OF F2

Rule	Confidence	
If <i>Moderate fees</i> then <i>Placement</i>	8/8	100 %
If <i>Moderate fees</i> then <i>Good result</i>	8/9	89%
If <i>placement</i> then <i>Moderate fees</i>	8/8	100 %
If <i>Placement</i> then <i>Good result</i>	8/9	89%
If <i>Good result</i> then <i>Moderate fees</i>	9/8	112.5%
If <i>Good result</i> then <i>Placement</i>	9/8	112.5%
If <i>Placement</i> then reference	8/3	267%
If <i>Placement</i> , then <i>state of the art lab</i>	8/8	100%
If <i>Placement</i> then <i>Library</i>	8/4	200 %
If <i>State of the art lab</i> then <i>Placement</i>	8/8	100%
If <i>distance</i> then <i>Transport facility</i>	2/4	50 %
If <i>distance</i> then <i>hostel facility</i>	2/3	67 %
If <i>Good result</i> then <i>State of the art lab</i>	9/8	112.5%
If <i>State of the art lab</i> then <i>good result</i>	8/9	89%

From Table 2 and Table 3 we can find all rules in this case. They are shown in Table 4 below.

TABLE 4 FINAL RULES

Rule	Confidence	
If Placement, Moderate fees then good result	6/9	67%
If Moderate fees, Experienced faculty then Placement	5/8	62.5%
If Moderate fees, state of the art lab then good management	1/2	50%
If Experienced faculty, state of the art lab then good result	2/9	22%
if Moderate fees, Placement Then good management	6/2	300%

III.CONCLUSION

We have described an algorithmic approach for discovering a dynamic pattern to improve the number of potentials is fast in various ways. It is exclusively useful for very large databases. An important contribution of our approach is that it drastically reduces the I/O overhead. This feature may prove useful for many real-life database mining scenarios where the data is most often a centralized resource shared by many user groups, and may even have to support on-line transactions. In near future we have a proposal to analyze the disk I/O and CPU overhead. In addition the new algorithm has an excellent scale-up property. The main drawback of this algorithm is estimate the number of partitions for the given available memory. This can be rectified in near future. We can extend this work by parallelizing the algorithm for a shared multiprocessor machine.

REFERENCES

- [1] S.ThabasuKannan, "Optimized mining of Very Large database via ClusteredIndexing Method", InternationalJournal of Intelligent Optimization Modeling, ISBN : 81-8424-104-6, Allied Publishers (p) Ltd, pages: 307 -318, 2009.
- [2] R. Agrawal and R. Srikant. Fast algorithms for mining association rules in largedatabases. In Proceedings of the 20thInternational Conference on Very LargeData Bases, Santiago, Chile, August 29-September 1 1994.
- [3] M. Houtsma and A. Swami. Set-oriented mining of association rules. InProceedings of the International Conference on Data Engineering, Taipei, Taiwan, March 1995.
- [4] A. Savasere, E. Omiecinski, and S. Navathe. An efficient algorithm for miningassociation rules in large databases. Technical Report GIT-CC-95- 04, Georgia.Institute of Technology, Atlanta. GA 30332, January 1995.
- [5] S.ThabasuKannan, "Knowledge Based Query processing of VLDB via ClusteredIndexing method", Proceedings of the International conference on GlobalManufacturing and Innovations with university of Massachusetts, Dartmouth –USA, collaboration with International journal of Operations Research, page: 155 - 158
- [6] K.C.C Chan, A.K.C. Wong and D.K.Y. Chiu, "Learning sequential patterns forprobabilistic inductive prediction," IEEE Trans. Systems, Man and Cybernetics, vol. 24, no. 10, pp. 1532-1547, 1994.
- [7] S.ThabasuKannan, "Discovering a pattern for effective utilization of Large ScaleDatabase via Clustered Indexing method", In proceedings of International Conference onKnowledge management and Information, organized by IADI Society, at Barcelona,Spain, pp. 518-525,2011.
- [8] SudiptoGuha, Rajeev Rastogi, and Kyuseok Shim. CURE: An efficient clusteringalgorithm for large databases. In ACM SIGMOD International Conference onManagement of Data, 1998.
- [9] Tian Zhang, Raghu Ramakrishnan, and MironLivny. BIRCH: An efficient dataclustering method for very large databases. In ACM SIGMOD InternationalConference on Management of Data, 1996.
- [10] Haisun Wang, Wei Wang, Jiong Yang, and Philip S. Yu. Clustering by patternsimilarity in large datasets. In ACM SIGMOD International Conference onManagement of Data, 2002.
- [11] S.ThabasuKannan, "An algorithmic approach for a simple prototype of business system to get customer satisfaction on CRM", International Journal of Business Review ISBN - 2249:5444, Vol 4, Issue 3, 2013.