

## Fast Cluster-learning with Prior Probability from Big Dataset

Tengyue Li, Simon Fong

Department of Computer and Information Science  
University of Macau, Macau SAR, China  
e-mail: mb75436@umac.mo, ccfung@umac.mo

Joao Alexandre Lobo Marques

School of Business  
University of Saint Joseph, Macau SAR, China  
e-mail: alexandre.lobo@usj.edu.mo

Raymond K. Wong

School of Computer Science and Engineering  
University of New South Wales, Sydney, NSW, Australia  
e-mail: wong@cse.unsw.edu.au

**Abstract**—Association Rule Mining by Apriori method has been one of the popular data mining techniques for decades, where knowledge in the form of item-association rules is harvested from a dataset. The quality of item-association rules nevertheless depends on the concentration of frequent items from the input dataset. When the dataset becomes large, the items are scattered far apart. It is known from previous literature that clustering helps produce some data groups which are concentrated with frequent items. Among all the data clusters generated by a clustering algorithm, there must be one or more clusters which contain suitable and frequent items. In turn, the association rules that are mined from such clusters would be assured of better qualities in terms of high confidence than those mined from the whole dataset. However, it is not known in advance which cluster is the suitable one until all the clusters are tried by association rule mining. It is time consuming if they were to be tested by brute-force. In this paper, a statistical property called prior probability is investigated with respect to selecting the best out of many clusters by a clustering algorithm as a pre-processing step before association rule mining. Experiment results indicate that there is correlation between prior probability of the best cluster and the relatively high quality of association rules generated from that cluster. The results are significant as it is possible to know which cluster should be best used for association rule mining instead of testing them all out exhaustively.

**Keywords**—Association Rule Mining; Clustering; Preprocessing; Prior Probability

### I. INTRODUCTION

Finding association patterns is becoming ever important analytics in the era of big data. Data that are accumulated from various sources, archived into a huge repository is a common practice nowadays due to the hype of big data. This trendy IT investment, is ranging from IoT oriented applications to smart city applications. Organizations that adopted big data processing technology store up the data hoping that insights or business intelligence could be harvested from that one day. Other than predictive

modelling, association rule mining is a popular approach in extracting association patterns among the numerous items or attribute values which characterize Tetra or Petra bytes of big data sitting in some data archive. Association patterns which are more technically known as association rules are statement-like kind of output that show the linkage between two or more attribute values. Users may find such association rules useful as they tell about which items or attribute values that often occur together by high frequency. Some example such as bank data would reveal how certain attributes of customers mostly relate or associate to certain purchases of financial products. Marketing personnel hence can find clues of what factors would lead to a customer buying a certain financial product or not buying a financial product. Supermarket, for another example is a typical scenario where association rules are usually generated. The association rules from supermarket products tell about which items are most frequently purchased together.

The association rule mining algorithms, such as Apriori method [1] for instance and other successors were invented in the late 90's. They were not designed to work with big data in mind in those days. Datasets from which association rule mining works on were relatively much smaller than the massive volume of big data that are emerging today. Apriori algorithm, being a classical association rule mining algorithm is known to have a drawback in scalability. The time performance in big O notation would have it scale up exponentially to the dimensions or the number of attributes of the data. So does the data volume of the big data which escalates the timing performance of association rule mining algorithms.

In order to tackle the challenges of association rule mining over big data or big dataset where the data are usually presented as a two-dimensional data matrix, some modifications are needed either at the algorithmic level or at the data processing level. The former implies that some new association rule mining algorithms need to be innovated or extended. The latter is relatively simpler because the existing algorithms could still be used, only the processing

or pre-processing stage of the knowledge discovery methodology is required to be improvised.

Earlier on, the authors advocated an alternative pre-processing method [2] which aims at improving the quality of association rules. It first segments the full dataset into clusters prior to running the association rule mining algorithm on. The underlying concept is to mine the association rules from subsets of the data which are supposedly more concentrated in data that occur frequently together. This is because clustering helps pulling data that are similar to each other together, reducing the distance space between the scattered data throughout the whole dataset. Therefore, it is easier to mine out the association rules from the frequently occurred data as they are close to each other. By this simple idea, the overall association rules quality improves. But the problem is, in advance there are no way to know which is most appropriate to be mined by association rule mining algorithm. Although the cluster size is smaller than the original dataset size, association rule mining algorithm such as Apriori still takes time. Furthermore, multiple clusters are often generated as a result of clustering analysis. Since it is not known in advance which cluster gives the best results, all the clusters are to be mined exhaustively by brute-force. This trying-all-out process is indeed time consuming.

In this paper, an indicator is proposed to be used as a hint for selecting a cluster out of all, as a candidate for the subsequent association rule mining. Prior probability is investigated in this case as a potential indicator in selecting the right cluster. The reminder of the paper is organized as follow: Section 2 reviews some prior arts on improving the quality of associate rules. Section 3 presents the methodology of the clustering-then-association-rule-mining. Evaluation experiment is documented, and the results are shown in Section 4. Section 5 concludes the paper.

## II. RELATED WORK

In the literature it is seen that many researchers worked on improving the quality of the association rules. In [4] a combination of association rule mining and means of rough sets when used together helps improve the overall quality of association rules. It was shown that identical rules are generated despite their different approaches because they are based on the same principle. Then it was shown that [5] using preprocessing works for improving the quality of association rules. In that study [6], a right choice of discretization method in preprocessing numeric data help improve the subsequent associate rule mining. The idea of preprocessing was taken further by this group of researchers in [7]. Clustering was first used here by detecting community in potential clusters. It was known that the main drawback of association rule mining is that many uninteresting rules were generated. Clustering algorithms were tried on the database for downsizing the dataset before association rules are extracted. It was argued in the paper that by only relying on the distance measures which are

known as similarity would have limited effect. Additionally, community detection algorithms are used. Several community detection algorithms were tested over two clustering methods. They showed that community detection algorithms perform well in this regard. The work in [8] shows there exists a strong relation between the quality of the dataset and the association rules. High-quality datasets have a very positive impact over the quality of the discovered association rules, and vice versa. Hence, integrate data quality measures for effective and quality-aware association rule mining and a cost-based probabilistic model for selecting legitimately interesting rules is proposed in [8]. Over most of the previous works reviewed, a relatively simple approach [9] is to use K-means clustering method to divide a large dataset into subsets. Then association rules for top 5 assignees are extracted from each cluster. The results show that essential rules are less than redundant rules in every cluster are generated. The limitation of this approach is that the user must specify a k value for partitioning a certain number of clusters. Different k value would lead to different results. An ideal k value is not known in advance unless exhaustive trial tests are run which of course is time consuming. Furthermore, when many clusters or subsets are generated, it is not known which the best is to use. All the works progressed previously are limited to some manual tuning or intervention in finding the right cluster out of many for subsequent association rule mining.

## III. OUR PROPOSED METHODOLOGY

In our earlier work [2], a cluster-then-associate workflow was proposed, and the experiment results show that the association rules are of better quality than those extracted directly from the original dataset. As shown in Figure 1, a direct association rule mining from the original dataset is at 1(a). Figure 1(b) depicts the workflow that was presented in [2]. In this paper, an additional mechanism is proposed to be there between the outputs of the clustering and the inputs to the association rule mining, as in Figure 1(c). The extra step is for choosing the most appropriate cluster among however many the resultant clusters or subsets there are at the end of the clustering step.

Traditional association rule mining loads the whole dataset for generating association rules. In Apriori algorithm, different itemsets are enumerated iteratively by scaling up the cardinality of the itemset length from zero to the maximum number of items that are purchased (or occurred) in a single transaction. Along the itemsets discovery process, itemset candidates are eliminated when they do not meet the minimum confidence or minimum support requirement. When the dataset is large, often brute-force exhaustive approach of keeping and eliminating the itemset candidates involves repeatedly scanning the whole dataset. To get around this time-consuming process, as indicated by earlier works in Section II, the rule generated process is controlled by harvesting only a certain number of

rules which usually are the top ones (e.g. top 20). Alternatively, a subset of dataset is used for the association rule mining instead of the whole one. In this case, out of several clusters or subsets of the original dataset that was partitioned by some clustering algorithm, there should be one or few out of all clusters should be used for incomplete association rule mining. See Fig. 1(b).

Fig. 1(c) shows a smarter way than that of Fig. 1(b) to use some indicator as clue in picking an appropriate cluster

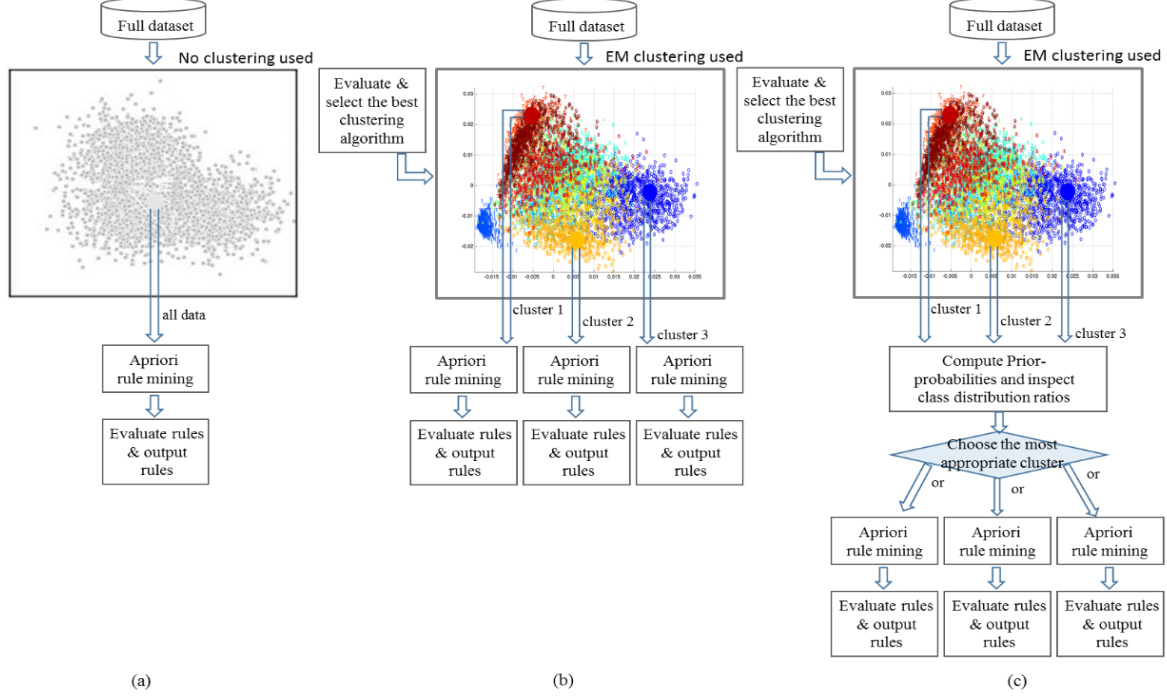


Figure 1. Cluster-then-associate workflow; (a) no clustering, (b) clustering only, and (c) clustering with cluster selection

Hence it is no longer necessary to try out every cluster, because a user can base on the Prior Probability (PP) to decide which cluster or a short-listed set of clusters to be selected for association rule mining. This process is therefore speeded up. Due to the speed saving, this association rule mining strategy is called Fast Cluster Learning (FCL). The meaning is to learn or to build a model quickly from a cluster in lieu of the full dataset. Without the need of using the full dataset, which is supported by an appropriate clustering algorithm, the subsequent model learning process can be done much more quickly especially for classical associate rule mining algorithm such as Apriori. It is known that Apriori when working with smaller dataset is fast, vice versa it may take a very long time when the dataset is large. In the era of big data, dataset which is subject to building or learning a model can be very large. FCL divides up the large dataset into small fragments, not randomly, but according to the similar data which shall be grouped together. Then a fast statistic indicator value is computed, PP it is in this case. With the PP value in place for each cluster, the cluster candidates can be sorted in descending order. The cluster candidate that has the highest PP value is suggested to be

out of the rest. The qualities of the clusters vary as there is no rule-of-thumb exists in either assuring the clusters are made all suitable for association rule mining or knowing in advance which cluster is most suitable for highest quality rules to be mined. As a solution, the proposed process in Fig. 1(c) advocates using a statistical property called Prior Probability to be used.

chosen for immediate association rule mining. Often it may yield the best quality rules. However, it is still being tested on the consistency. So, the user may be advised to shortlist several clusters from the top of the list, for rule generation; then rules are generated from each of the shortlisted clusters (hopefully not too many), for inclusion into the final rule result list.

In the methodology, parameter free clustering algorithm namely Expectation Maximization (EM) algorithm is suggested to be used. EM is known to be one of the most powerful clustering algorithms which are without needing the user to specify how many clusters that the algorithm needs to generate. It will automatically find the optimal number of clusters by its maximization mechanism. The mechanism is the core design by probabilistic clustering. The design of probabilistic clustering or learning takes on this assumption –data items should have certain affinity in terms of probability belonging to all the clusters; the data items shall not exclusively belong to any particular cluster. By this principle, the assumption basically states that certain amount of evidence is needed to decide on how clustering should be done, even by probabilities. Some special properties are

assumed – probability distributions are used to describe each cluster. These probability distributions are computed over the attributes of the data items which are the cluster members. The probability distributions therefore can be seen or used as deciding factors on the choice of cluster membership for each data item. Consequently, an important concept called Prior Probability (PP) is used to represent such probability distribution; not all clusters have identical PP because the clusters have different probability distributions with respect to different data items.

Let  $pp_j$  be the PP of the cluster  $C_j$  where  $C_j$  has some mean value  $\mu_j$  and standard deviation  $\sigma_j$ , in such a way that the cluster  $C_j$  can be described statistically in some shape of a Normal distribution. For brevity, the data items only have a single numeric attribute whose values distributed Normally over all the  $k$  clusters,  $C_1, C_2 \dots C_j \dots C_k$  where  $j \in [1, k]$ . In a very primitive state of probabilistic clustering, the following are the parameters of the model that need to be calculated:

$$Model\_Param = Model(pp_j, \mu_j, \sigma_j, \text{ where } j \in [1, k]) \quad (1)$$

The set of data items are iteratively evolving from previously belonged clusters, subject to be reassigned to possibly different clusters. At the beginning the *Model\_Param* values are set randomly; based on the given the current parameter values, the probability of the cluster membership for each data item is computed; the *Model\_Param* values are re-evaluated and updated using the calculated probabilities which are shown to have yielded better cluster memberships. This iterative process repeats until the *Model\_Param* appears to converge where no further improvement on the cluster membership can be observed. As such the EM clustering algorithm consists of two key phases which alternate until convergence:

Phase (1) Expectation step, where the cluster  $pp_j$  is computed in each  $j^{\text{th}}$  iteration based on the model parameters values. The probability of each data item  $x_i$  where  $i \in [1, n]$  from the dataset is computed for finding out the best cluster membership for the cluster  $c_j$  where  $j \in [1, k]$  using the Equation (2):

$$Pr(E) \forall_{i,j} = \varepsilon_{i,j} = pp_j \times Pr(x_i | c_j) \quad (2)$$

where  $Pr(x_i | c_j)$  is calculated from the Normal distribution of the  $j^{\text{th}}$  cluster where  $j \in [1, k]$  given that the current model is configured with the *Model\_Param*  $\in (x_{*j}, \mu_j, \sigma_j)$ ; it is noted that in this situation all the model parameter values unknown directly, but they could be inferred via PP.

Phase (2) Maximization step, calculates the most suitable *Model\_Param* values in order to maximize the likelihood of models that are guessed to match the current positions or memberships of the data items. The model parameter values  $Model(pp_{j,j}, \mu_{i,j}, \sigma_{i,j}, \text{ where } i \in [1, n] \text{ and } j \in [1, k])$  in this step are re-calculated using the following formula:

$$\text{Prior probability: } pp_{i,j} = \sum_i \frac{\varepsilon_{i,j}}{n} \quad (3)$$

$$\text{Mean: } \mu_{i,j} = \frac{\sum_i x_i \times \varepsilon_{i,j}}{\sum_i \varepsilon_{i,j}} \quad (4)$$

$$\text{Standard deviation: } \sigma_{i,j} = \sqrt{\frac{\sum_i |x_i - \mu_j|^2 \times \varepsilon_{i,j}}{\sum_i \varepsilon_{i,j}}} \quad (5)$$

These dual steps repeatedly increase the log-likelihood of all the clusters until there is no more significant refinement according to the Eqn.(6).

$$\log Pr(x) = \log \sum_j (Pr(x | c_j) \times pp_j) \quad (6)$$

Usually the overall quality of the clusters which is represented by log-likelihood will rise sharply at the beginning over some initial iterations. The EM will then converge to an equilibrium state with a stable log-likelihood value. For picking the most appropriate clusters out of all, log-likelihood cannot be used here because it is an indicator that sums the qualities of all the clusters. However individually, the prior probability,  $pp_j$  indicates the fitness of each cluster. For this reason, in our methodology  $pp_j$  is used as a measure in choosing the priority cluster for subsequent associate rule mining.

Furthermore, in our fast cluster learning methodology, EM is chosen because it is guaranteed to converge to highest possible log-likelihood fitness. The only limitation is the time consumption which can be significantly large. It is because there exist local maximum and global maximum, similar to that of k-means [10]. For prevention of falling into local maximum, the EM algorithm is programmed to run several times for obtaining some chances of reaching the global maximum as each time they start with different random orientation for the initial clusters and model parameters, guessing what their suitable parameter values are. It is the compromise between using a heavy algorithm that takes considerable amount of time to run and achieving the best clusters most of the time without specifying the  $k$  parameter value.

#### IV. EXPERIMENT AND RESULTS

For validating our proposed methodology, three empirical datasets and two artificially generated datasets are used. The objective is to investigate the associate rule mining performance over the fast cluster learning method. The comparison baseline is standard Apriori applied over the full dataset. In the cases of fast cluster learning, EM is used to partition the full dataset into clusters, and prior probability is used as the quality indicator in selecting only the most suitable cluster for rule mining. The three datasets in use are: bank data which contain customers' attributes such as simple demographic information and income levels, and class target of whether this customer has purchased an investment product; homicide data that breaks down each home violence case with characteristics of how it happened; and lung cancer data that show the simple backgrounds, lifestyles and eating diets of patients who are confirmed to have lung cancers and

those who are negative. The datasets are available for public download at the UCI machine learning data repository.

Some assumptions and pre-processing are done: in the setting of the Apriori algorithm, the algorithm parameters are: *classification-association=true, minMetric=0.7, numRules=20*; numeric attributes are converted to nominal attributes by binning the numbers in regular intervals. The performance in comparison is the averaged confidence values over the top 20 rules extracted. The number of interesting rules is also used in the case when the mode of classification-by-association is turned on – that is when the user is looking for rules that associates with the target class, e.g. *class=yes* in medical data. The results of the associate rule mining of the three datasets using the normal mining method without clustering and that was clustered in advance (as in Fig. 1(b)), are shown in Figs. 2, 3 and 4 respectively.

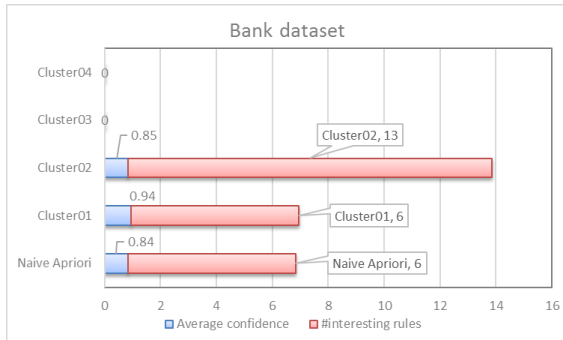


Figure 2. Fast cluster learning vs naive Apriori over Bank dataset

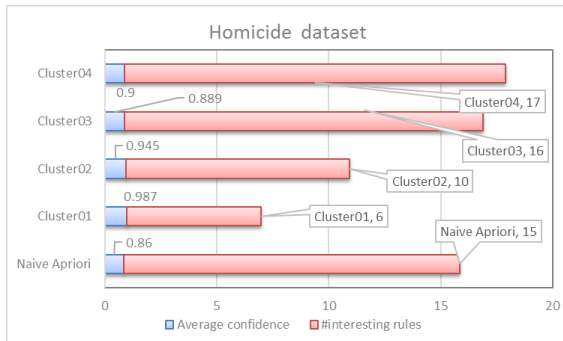


Figure 3. Fast cluster learning vs naive Apriori over Homicide dataset

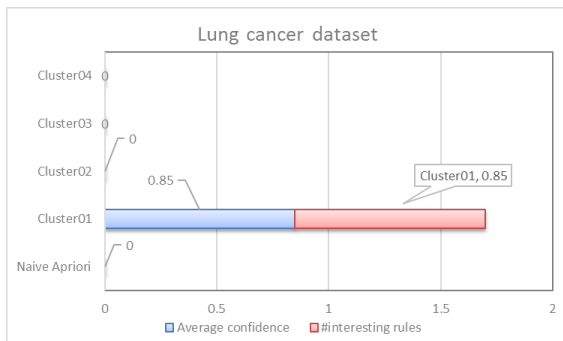


Figure 4. Fast cluster learning vs naive Apriori over Lung cancer dataset

It can be observed that over the three cases of using different empirical datasets, fast cluster learning which learns the association rules from most of the clustered data have superior results (in average confidence and number of interesting rules) to the Naïve Apriori results. There is always at least one set of results which are harvested from certain clusters that are better than the original results by Apriori without clustering. The next research question naturally would be ‘which cluster to select?’. Since it is not known until all are to be tried for finding the ground truth, in this paper we advocate using Prior Probability as a clue in the cluster selection.

First of all, the hypothesis that the Prior Probability (PP) is related to the quality of the resultant association rules. In this case, the quality is defined as average confidence of the top rules. Hence, a scatter plot is charted showing the relation between PP and average confidences of the rules over the three cases of empirical datasets. The results are charted in Figure 5.

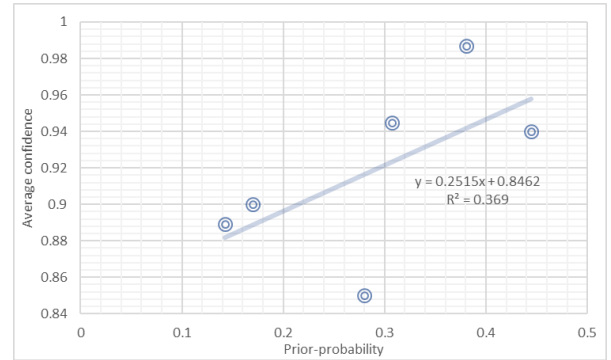


Figure 5. Correlation between prior probability and average confidence by the three empirical datasets

To further verify whether there are relations between PP and average confidence, artificially generated data are used over which we can control the variation of data sizes. Two artificial datasets are generated for this purpose, one is by increasing the data rows from 100 to 10 million instances, with a step enlargement of 10 times at each increment. The number of features or attributes is kept at 10. The other artificially generated dataset increases its dimension by expanding the attributes from 10 to 50 with a step size of 10. The purposes of using these two artificial datasets are: 1) observe whether there exists a correlation between the PP and the quality of rules (in this case, minimum support is used for generating top 20 rules); the other purpose is to investigate the scalability of our fast cluster learning methodology by PP. Particularly, two main parts of time consumptions are apparent and being observed here. They are the time consumed mostly by the clustering and by the Apriori rule mining. Figure 6 and Figure 7 show the relation between PP and the quality of association rules by using artificial datasets that are scaled in data instance expansion and dimension or attributes increase respectively.

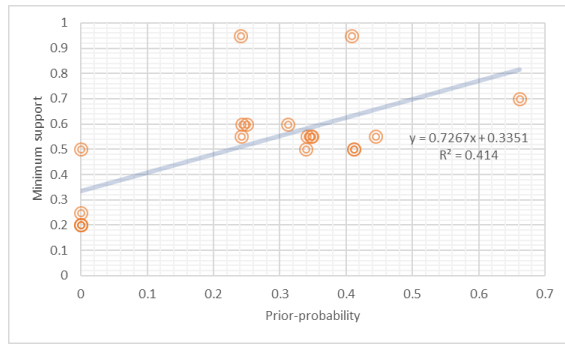


Figure 6. Correlation between Prior Probability and Minimum Support by the artificially generated data with escalation of data rows

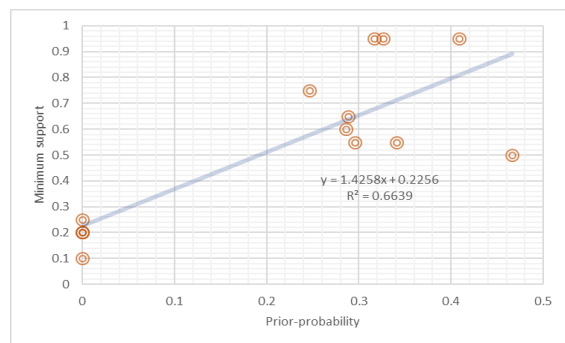


Figure 7. Correlation between Prior Probability and Minimum Support by the artificially generated data with escalation of data attributes

In Figures 5, 6 and 7, Pearson correlation coefficients are computed, and a linear correlation line is added. The Pearson correlation coefficient are 0.369, 0.414 and 0.6639. It shows that there are somewhat correlations between PP and the quality of rules. However, in general, for the case of empirical data, the correlation is not too strong probably because the underlying patterns of the empirical datasets are nonlinear and random, relatively. For the artificially generated datasets that are outputted by a standard random generator, they show stronger correlation. Particularly, when the dataset is varying in the dimensions or attributes, the correlation that that of varying by data rows. This implies high dimensional data with attribute values scatter far apart in the data space, clustering has more significant effects in bringing the useful data items closer into a cluster. Hence Apriori that applies on the concentrated cluster yields better results. The time performance in the scalability test is carried out, comparing k-means, EM and Apriori method. It is to observe whether it is feasible to adopt such methods in the fast cluster learning methodology, especially when the data row size and the dimension escalate up greatly in the era of big data which is not uncommon. Figures 8 and 9 are the time consumption charts for showing the time required by different data mining mechanisms, over the artificial dataset that grows in data volume and dimensions respectively.

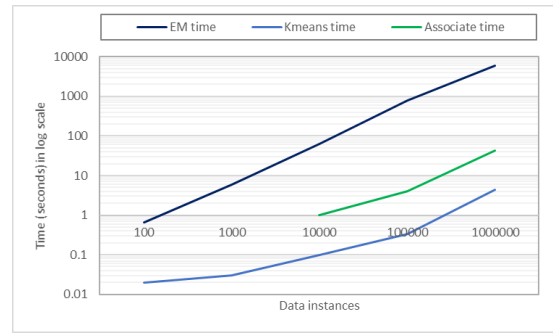


Figure 8. Time consumption by various processing methods in Fast Cluster Learning over the artificially generated data with escalation of data rows

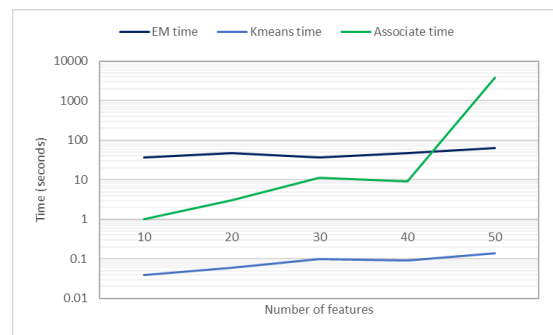


Figure 9. Time consumption by various processing methods in Fast Cluster Learning over the artificially generated data with escalation of data attributes

In Figure 8, the three processing methods all scale up somewhat linearly. That means the time increase is proportional to the increase of data volume. EM is increasing at a rate higher than the other two, possibly because EM is an optimization algorithm which iteratively improves the model parameter values; more data instances the longer time EM needs to converge. In Figure 9 there are interesting phenomena; the EM and k-means are rather flat. Especially it is the EM which scales up very well. The other interesting phenomenon is the exponential time increase by Apriori association rule mining method which is well-known in research community. When the data attributes increase, Apriori will increase explosively at certain threshold. There is a cross point somewhere near the attribute size equals to 42. At that point, both EM and Apriori take about the same amount of time in processing. However, when the attribute size goes beyond 42, Apriori time will shoot up greatly and EM remains as a steady time processing overhead. The experiment shows that EM is not sensitive to attribute increase but data volume increases, and vice-versa for Apriori.

## V. CONCLUSIONS

In this paper, a new data mining methodology called Fast Cluster Learning is proposed. It is designed for enhancing

the quality of rules by association rule mining such as Apriori method. Previously it was shown that by using clustering to first partition the data in clusters help enhance the associate rules at the output. But the problem was it wasn't easy to know which cluster to use given there are multiple clusters to be generated after EM is run. In this paper, a simple quality indicator called Prior Probability (PP) is proposed to use for quickly identifying a cluster that would be useful for subsequent rule mining. Experiments are conducted using both empirical datasets and artificially generated datasets. All the cases of datasets point to a positive correlation of different strengths between PP and the quality of the rules.

#### ACKNOWLEDGMENT

The authors are thankful to the financial support from the research grants, MYRG2016-00069, and FDCT/126/2014/A3, offered by FDCT of Macau SAR government.

#### REFERENCES

- [1] Rakesh Agrawal and Ramakrishnan Srikant Fast algorithms for mining association rules. Proceedings of the 20th International Conference on Very Large Data Bases, VLDB, pages 487-499, Santiago, Chile, September 1994.
- [2] Simon Fong, Robert P. Biuk-Aghai, Scarlet Tin, Visual clustering-based apriori ARM methodology for obtaining quality association rules, Proceedings of the 10th International Symposium on Visual Information Communication and Interaction, Bangkok, Thailand, August 14-16, 2017, pp.69-70
- [3] Paul W. Purdom, Dirk Van Gucht, and Dennis P. Groth, Average-Case Performance of the Apriori Algorithm, SIAM J, 33(5), pp.1223-1260
- [4] Daniel Delic, Hans-J. Lenz, and Mattis Neiling, Improving the Quality of Association Rule Mining by Means of Rough Sets, Soft Methods in Probability, Statistics and Data Analysis. Advances in Intelligent and Soft Computing, vol 16. Physica, Heidelberg, pp.281-288
- [5] Maria N. Moreno, Saddys Segrera, Vivian F. López and M. José Polo, Improving the Quality of Association Rules by Preprocessing Numerical Data, II Congreso Español de Informática, pp.223-230
- [6] Liu, H., Hussain, F., Tan, C.L., Dash, M. Discretization: An Enabling Technique. Data Mining and Knowledge Discovery, 6, 2002, pp.393-423
- [7] Renan de Padua, Exuperio Ledo Silva Junior and Laes Pessine do Carmo, Veronica Oliveira de Carvalho, Solange Oliveira Rezende, Preprocessing data sets for association rules using community detection and clustering: a comparative study, XIII Encontro Nacional de Inteligência Artificial e Computacional, 2016, pp.553-564
- [8] Laure Berti-Equille, Quality-Aware Association Rule Mining, PAKDD 2006: Advances in Knowledge Discovery and Data Mining, pp.440-449
- [9] Meera Sharma, Abhishek Tandon, Madhu Kumari and V. B. Singh, Reduction of Redundant Rules in Association Rule Mining-Based Bug Assignment, International Journal of Reliability, Quality and Safety Engineering. Vol. 24, No. 06, 1740005 (2017)
- [10] Simon Fong, Suash Deb, Xin-She Yang, and Yan Zhuang, Towards Enhancement of Performance of K-Means Clustering Using Nature-Inspired Optimization Algorithms, The Scientific World Journal, vol. 2014, Article ID 564829, 16 pages, 2014