



Efficient streaming data association rule mining

Amr Aly Abd Elaty^{a*}, Rashed Salem^b, Hatem Abd Elkader^c

^a Faculty of Computer and Information, Information Systems Department, Menoufia University, Menoufia, Egypt,

^b Faculty of Computer and Information, Information Systems Department, Menoufia University, Menoufia, Egypt,

^c Faculty of Computer and Information, Information Systems Department, Menoufia University, Menoufia, Egypt.

KEYWORDS

Data Mining,
Association Rules,
Streaming Data,
FUPP Tree.

ABSTRACT

Recently, number of applications including social networks, stock market trading and sensor network devices generate a massive amount of data in the streaming form. Streaming data have characteristics different from static data, such as streaming data arrives continuously at high speed with huge amount. Mining and discovering information from these data is a non-trivial issue. Most of traditional algorithms have limitations to deal with streaming data, so there are new issues raised and need to be taken into consideration while developing techniques for mining association rules from such data. In this paper, a technique to mine an association rules from streaming data efficiently is proposed. The proposed technique develops a tree structure called Fast Update Frequent Pattern Tree (FUPP-Tree) that reduce the number of traversing between tree nodes in both inserting a new transaction and extracting an association rules between items. Also, to avoid congestion during inserting incoming streaming data to FUPP-Tree, a sliding window approach is used to divide incoming data equally to all available windows. The complexity and the performance of this technique are investigated, and a dataset of storehouse is used to test the proposed technique and measure its efficiency. The efficiency of the proposed technique is compared with other most related algorithms.

1. Introduction

Over the past years, some applications, for example, social networks, stock market trading and sensor network devices need to process data as they are generated, in other words, as they stream. These types of applications are called streaming data applications. The term of streaming data refers to data that is generated continuously with unbounded size and arrives in high speed, as opposed to static data. Stream data mining is a process of extracting knowledge from rapid and continuous incoming data, it portrays the next era of data mining systems that will enable the intelligent and time-critical information requirements of portable users and will ease "anytime, anywhere" data mining [1]. Association rule mining technique is aiming to discover a frequent patterns, correlations or associations from a given dataset. It is termed as "market basket dataset", where each attribute is termed as an item and the frequencies of different itemsets are transformed in the form of if-then rules based on support-confidence framework and then the relationships between seemingly unrelated data in given dataset

can be found out [2]. Streaming data applications require association rule mining to discover the major associations among items. Some of related algorithms are discussed, there are limitations for these algorithms to deal with streaming data and mining an association rules in high performance and in acceptance time. The proposed technique tackles the discussed algorithms limitations and uses a sliding window to build an enhanced FP-Tree called FUPP-Tree. Also, the proposed technique allows fault tolerance layer that aims to save the resulted tree after a period of time to retrieve it in the case of any error or damage to the system rather than rebuild the whole tree again from scratch. Saved tree can be used to apply user query for generating the association rules using both minimum support and minimum confidence thresholds. The proposed technique is applied for a dataset and the experimental results proved that it is efficient compared to other similar algorithms. In this paper, the related works are discussed declaring the limitations of each algorithm in section II, then the proposed technique is developed based on the discussed limitations in section III. The algorithm complexity is discussed in section IV. Moreover, the

* Corresponding author. Tel.: 0201113531165.

E-mail address: amrally_267@yahoo.com

Peer review under responsibility of Dr. Rashed Salem.

experimental results and querying the proposed tree are discussed in section V. Finally, the conclusions and future trends are discussed in section VI.

2. Related Works

The main issue in mining the frequent itemsets in streaming data is to determine the frequency of them at a suitable rate that is well-suited with the speed at which the transactions are provided. This objective needs algorithms featured with in-memory data structures and a minimal dataset scan. According to [2-3], the approaches of stream mining can be categorized into four main classes: bottom-up, top-down, landmark and sliding-window based mining.

2.1. Bottom-Up approach

In this approach, the individual transactions itemsets of the given dataset are specified in detail firstly. Then these itemsets are linked together to form larger sub-transactions in many levels and so on until complete top-level sub-transactions are formed. The common association rules algorithms in this approach are Apriori Algorithm and Partitioning Approach.

2.1.1. Apriori algorithm

Apriori algorithm is one of the most common association rule algorithms [4]. It can be used to derive all possible frequent itemsets from a given dataset and generate association rules subject to support and confidence values that are not less than a predefined minimum support value and minimum confidence value. In spite of the fact that it is powerful as the noisy data can't influence the result of the algorithm, but it has limitations when dealing with streaming data. The major limitation of this algorithm is the several scanning of the dataset when there is a new record is inserted. Also, it's a costly waste of time to generate a number of candidate sets with much frequent itemsets, so that it will be very slow and inefficient when memory capacity is limited.

2.1.2. Partitioning approach

It can be observed that during the frequent itemset generation, maximum time is consumed while reading the data from the disk. To execute faster, the dataset need to be loaded to the memory. But in most of the cases the dataset is too big to load into the memory. The partitioning approach uses the Apriori algorithm for memory resident data [2]. In this approach, the whole dataset is splitted into some smaller partitions, so that each partition is individually loaded in the memory. Then for each partition, a frequent itemsets are generated using the Apriori algorithm. After the generation of frequent itemsets for all the partitions, they are combined together, and redundancies are removed. Then for all the remaining itemsets the support is counted by reading the dataset again. This approach practically requires two scans of the whole dataset. This approach still suffers from limitations such as it isn't sensitive to noisy data. This approach also scans the dataset only twice. Moreover, in the final phase, joining of frequent itemsets of individual partitions results in a huge number of itemsets and hence consumes a significant amount of time.

2.2. Top-Down approach

In this approach, an overview of the given dataset is formulated, then breaking down to gain insight from given dataset. The most popular association rules algorithm in this approach is FP-tree algorithm.

- **Frequent-Pattern growth (FP-Growth) algorithm**

It is an efficient tree-based algorithm to discover the required association rules. This algorithm firstly scans the dataset to count frequencies of different items, then it reorders the items based on the frequency of each item in the decreasing order. By utilizing the frequency descending list, the dataset is compacted into a Frequent-Pattern tree, which keeps the information about the association of the transaction itemsets. Next, for each item starting with the highest support, a conditional pattern base is constructed and represented as its conditional FP-tree. The growth pattern is realized via the chain of the suffix pattern with the generated frequent patterns from the conditional FP-tree. After the construction of the FP-tree, for every frequent item one conditional FP tree is constructed. However, a major limitation of FP-growth is that, this algorithm needs to scan the given dataset twice [5]: First scan to get frequency of occurrence for each item, second scan to reorder the dataset transaction items according to the frequency of occurrence of each item.

2.3. Landmark approach

In landmark techniques, the itemsets of incoming transactions are calculated among a specific timestamp, the landmark, and the present. Therefore, in such landmark techniques, transactions are continuing in the frame of interest. Landmark techniques are based on a single pass support count of streaming data as well as on prefix tree-based pattern representation [3]. DSM-FI algorithm is a popular algorithm that is based on landmark approach.

- **Data Stream Mining for Frequent Itemsets (DSM-FI) algorithm**

In this algorithm, it constructs and maintains an in-memory prefix-tree based data structure summary, called summary frequent itemset forest (SFI-forest). A DSM-FI algorithm prunes infrequent itemsets from the current SFI-forest. Finally, the frequent itemsets from the current SFI-forest are generated [6]. The major limitation of this algorithm is that it needs more tree traversals for the frequency count, so that it consumes more time in both inserting and generating an association rules.

2.4. Sliding-window based mining

The major issues are escaping several scans because the streaming data come from one source or multiple sources in a high speed. So that this technique is based on the sliding window model, which entirely ignore old data and attention is focused on recent data, thus saving memory storage and simplifying the discovery of the distribution drift [7]. There are many algorithms that using sliding window-based mining such as Weighted Sliding Window (WSW) algorithm.

- **Weighted Sliding Window (WSW) algorithm**

This algorithm depends on the number of windows for mining, the window size and the given weight for each window, which are predefined. The incoming transactions are split into equal number of windows, and then compute the weight of every transaction in every window. Hence, the highest weight has been assigned to the most recent transaction. If the

Table 1 – Comparative analysis of frequent pattern algorithms.

Algorithm Name	Category	Advantages	Limitations
Apriori Algorithm	Bottom-Up	Can derive all possible frequent itemsets from a given dataset and generate association rules subject to minimum support and minimum confidence values	Multiple scanning of the dataset when there is a new record is inserted
Partitioning Approach	Bottom-Up	Generate frequent itemsets faster without burden on the memory	Not sensitive to noisy data
Frequent-Pattern growth (FP-Growth) algorithm	Top-Down	Find out frequencies of different itemsets, then order the itemsets descending into a compressed frequent pattern tree	Needs to scan the given dataset twice
Data Stream Mining For Frequent Itemsets (DSM-FI) Algorithm	Landmark	Compact tree structure has been designed to store the frequent patterns	It needs more tree traversals for the frequency count
Weighted Sliding Window (WSW) Algorithm	Sliding window based mining	A single pass algorithm was developed to discover the frequent itemsets	Weights of each window affected the mining results. So, user should specify the reasonable weight for each window

weighted support value count of a specific item is not less than the minimum weighted value, it is called as frequent itemset. There is a tradeoff between the window size and the execution time of WSW, since when the window size is small, the number of transactions involving frequent itemsets in every window is also small. The main limitation of this algorithm is that weights of each window influenced the results of the mining process, so that user should determine the well-suited weight for every window and adjust the weights values for different windows depending on the significance of the data. All discussed related works are summarized and compared in Table 1.

3. Efficient Association Rules Mining With A Fault Tolerance

In this section, a proposed technique is discussed. The main purpose of the proposed technique is improving a tree structure which can accommodate streaming data and change continuously, so that, a Fast-Updated Frequent Pattern Tree (FUFPTree) is proposed. Furthermore, a Sliding-Window technique is used to speed up preprocessing of incoming streaming data before sending it to FUFPTree in a parallelism form with a fault tolerance level. Finally, the association rules between items can be extracted easily from the built FUFPTree according to given parameters such as minimum support value, maximum support value and confidence value. All these issues are detailed in the following sections.

3.1. Streaming data preprocessing using sliding-window

With the exponential growth of streaming, an unprecedented amount of structured, semi-structured, and unstructured data is available. So that, data preprocessing is a major phase to solve incoming data problems before insertion in a tree such as inconsistencies, missing values and noise data to provide a high quality data to improve the performance of used algorithms to extract the association rules. To handle continuous data streams, a model based on sliding window is utilized for parallel preprocessing entry data. Typically, the incoming data will be split into equal chunks according to

the window slide size. Given a parallelism degree in the system, each window slide will act as a separate part and apply a set of data preprocessing operations such as data cleaning, data integration, data transformation...etc. A common sliding window technique called "Pane-based Partitioning" is used which based on the panes [9]. The main idea of this technique is to split overlapping windows internally into individual panes, over which sub-aggregates can be calculated whose results can be merged into the final aggregate. The panes technique has been originally introduced to minimize both computation cost and the space of sliding window by sub-aggregating computation. The number of sliding windows, sliding window size and sliding window pane size must be identified firstly before receiving data from sources. In "Pane-based Partitioning" technique each window acts as a shared-nothing cluster of commodity hardware in which each window is independent, self-sufficient according to the window size. There is no shared memory or disk storage, so that if there is any failure in any sliding window, other sliding windows will continue their work and not affected by the existing failure. A sliding window can become out of service due to a common failure called "Fail-Stop Failure".

- **Fail-stop failure**

This failure means that if the sliding window gets out of services during a computation, it requires an external impact to bring the sliding window back to working state again. For instance, a system administrator checks for the status of sliding window components and solve the problem or retires the broken sliding window and reconfigure the system such as reinitializing the number of sliding windows.

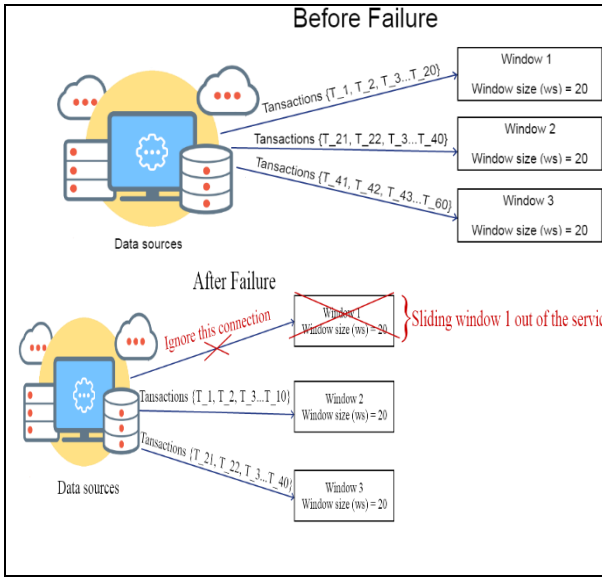


Fig 1 – Sliding window fault tolerance.

3.2. Sliding-window fault tolerance

Fault tolerance is the property that enables systems to remain operating correctly in the case of the failure of some of their portions due to one or more faults. As mentioned previously, the streaming data comes in continuously and with high speed rate, so that there is no ability to hold longer until the faults are processed from system administrator. In the statue of Fail-Stop Failure, the proposed solution is that the system will automatically disable and ignore the connection to the failure sliding window, so that all incoming transactions will be split directly to other working sliding windows without change the initial configuration of the sliding windows until the system admin interferes to resolve the issue with an appropriate solution. For instance, assume there is a system consist of n sliding window and each window can handle m transaction (i.e., window size = m), if one sliding window became out of service, that means there are $n-1$ working sliding window rather than n working sliding window and the transactions of the crashed sliding window will be redirect automatically to the other working sliding windows, see Fig 1.

3.3. Fast Updated Frequent Pattern Tree (FUFPTree)

The FUFPTree building algorithm is depend on the FP-tree algorithm. The connections between parent nodes and their children nodes are bi-directional linking that help to speed up the maintenance process such as reorder the tree elements rather than rebuild the tree from scratch. The FUFPTree structure depends on that the itemsets with the most frequent values will be in the top nodes as a descending order, so that, the value of bi-directional linking appears here when the itemsets frequencies values are changing and there is needing to reorder the tree rather than rebuild it [13]. Moreover, the frequent items are arranged in descending order and kept in the top nodes. Also, FUFPTree saves the last path that eases moving from and to nodes. For example, if last inserted transaction is {B, D, A, E}, then the saved path is {B:5, D:3, A:2, E:1} as shown in Fig 2. If there is a new transaction {B, D, A, E, C}, there is no need to return to Null node and start from scratch to search for the header node, i.e., {B}, the saved path can be used as its items are found in the new incoming transaction, so the support of each item in the current saved path will be incremented as {B:6, D:4, A:3, E:2}. Then after adding a new item node {C} with its frequency, now

the current saved path is updated to {B:6, D:4, A:3, E:2, C:1} as illustrated in Fig 3. If the new transaction starts with an item that doesn't match the first saved path item, the return to Null node and search for header node that match the first transaction item is inevitable. If new incoming transactions are inserted, the FUFPTree maintenance algorithm will manipulate them to keep up the FUFPTree. The top nodes of FUFPTree are refreshed whenever needed. If an originally item with large value turn out to be smaller, it is directly transferred from top level node to lower level node with respect to its frequency value and its parent and child nodes are then linked together. The FUFPTree can be updated, and the performance of the FUFPTree algorithm can be incredibly enhanced. The whole FUFPTree can then be re-built in a batch manner when an adequately extensive number of incoming transactions have been inserted.

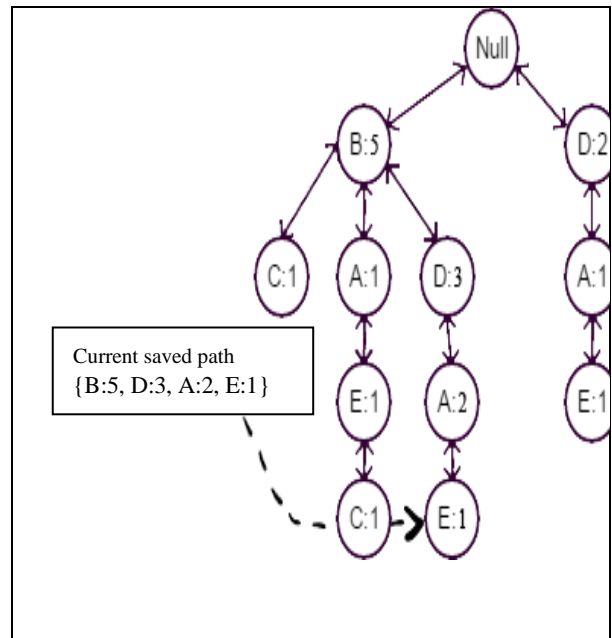


Fig 2 – FUFPTree before inserting {B, D, A, E, C} transaction.

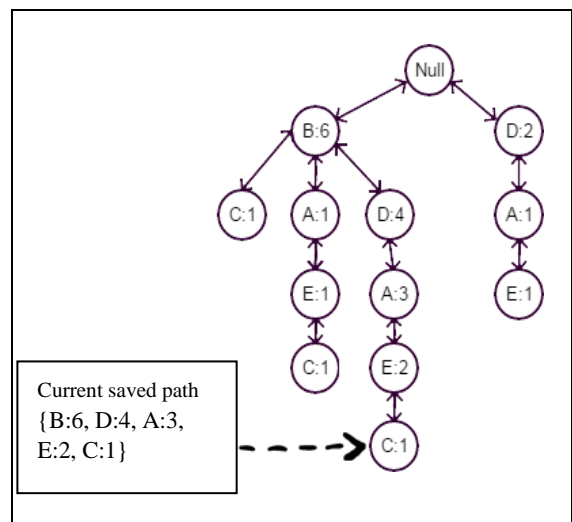


Fig 3 – FUFPTree after inserting {B, D, A, E, C} transaction.

3.4. Extracting association rules

Association rule is a data mining method for discovering interesting relations between different items. After constructing the FUPP-Tree, Algorithm A.1 is used to generate the entire set of correct frequent patterns from the current window [1]. Firstly, a minimum support threshold value must be identified to extract a frequent pattern list. In FUPP-Tree, the nodes are ordered in descending order, so that if the top node is not greater than or equal to the minimum support threshold value, this node and its children will be ignored from frequent pattern list. For example, from the same FUPP-Tree which shown in Fig 3, if the minimum support threshold value is 4, the resulted frequent pattern list will be {B, D}. After generating all possible frequent patterns, association rules are derived by a specific confidence threshold using Algorithm A.2.

```

Input: FUPP-Tree, Predefined minimum support threshold min_sup
Output: The frequent patterns list Lfp
Begin
1- For each node N in first level of FUPP-Tree
2-   If frequency(N) >= min_sup Then
3-     Extend the Node N
4-     For each item I in Node N
5-       If frequency(I) >= min_sup Then
6-         Add item I to Lfp , Then goto next level
7-       Else
8-         Break
9-       End If
10-    End For
11-  End If
12- End For
    
```

A.1 – Extract frequent patterns.

```

Input : Lfp list, Predefined confidence (Conf)
Output: The association rule list Lar
Begin
1- For each item Li in Lfp
2-   For each item Lj in Lfp
3-     Rule_val = Sup(Li U Lj)/Sup(Li)
4-     If (Rule_val >= Conf) Then
5-       Add Li and Lj to Lar list
6-     End If
7-   End For
8- End For
    
```

A.2. – Deduce association rules.

4. Complexity Of Algorithm

In this section, the algorithm complexity is analyzed. Runtime complexity is divided into two steps.

- Data preprocessing step which is depending on the number of sliding windows and the number of transactions per a sliding window (n), so that the complexity of sliding windows is O(n).
- Generating frequent patterns list from FUPP-Tree step which its complexity is O(log n).

5. Experimental Results

To evaluate the validity of the proposed technique, an experiment is conducted. A sample of a retail market basket dataset of a retail supermarket store is used. The number of transactions is 88163 and the number of unique stock keeping unit products is 16470 [10].

5.1. Sliding windows configuration

The main parameters for sliding windows are: number of windows (nw), window size (ws), number of panes per window (wp), size of pane per each pane in window (ps). All of these parameters are determined based on machine capabilities such as the machine operating system, hard drive type, hard drive capacity, processor and memory capabilities. In this experiment, the used machine specifications are as follow; Ubuntu 18.10 OS, solid state drive (SSD) with 100GB available capacity, processor is Intel core i5, 2.40 GHz and 4 GB RAM. Sliding windows configuration must be done quite ideally, because these sliding windows have a major role as they receive a raw data from sources and then apply a preprocessing function to be ready to be inserted into the FUPP-Tree. So that one of the best tools that is used during this experimental is called “Apache Kafka”. Kafka is an open source tool that is used for building streaming applications. This tool is scalable, fault-tolerant and allows executing operations fast. In this experiment, the parameters are configured with different values in proportion to the capability of the machine and deduce the execution time as shown in Table 2.

5.2. Building FUPP-Tree using pane-based partitioning VS. building FUPP-Tree using weighted sliding window

In this experiment the FUPP-Tree building time is measured by applying the proposed technique and applying weighted sliding window (WSW) algorithm. As mentioned previously, WSW algorithm needs to identify number of windows for mining, the size of the window as well as the weight for every window firstly. The major limitation is that the weights of every window influence the mining results. It requires very reasonable weight for each window and adjusts the weights values each time for each window, so that it takes more time to construct FUPP-Tree than using pane based partitioning approach. In pane-based partitioning approach, there is no need to adjust the configurations that are identified before starting, so that the FUPP-Tree is constructed faster as shown in Fig 4.

Table 2 – Total runtime with different sliding window configuration.

Number of windows (nw)	Window size (ws)	Number of panes per window (wp)	Size of pane per each pane in window (ps)	Total Runtime (sec)
10	8800	80	110	20.5
20	4400	50	88	12.9
30	2930	10	293	8.7
40	2200	20	110	7.4
50	1760	10	176	4.2

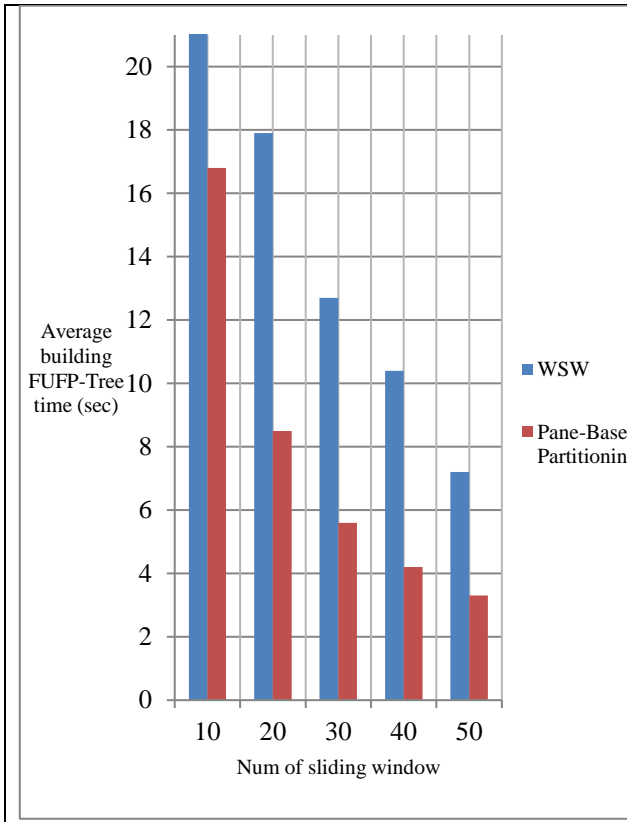


Fig. 4 - Average time for building FUFP-Tree.

Table 3 – Top 5 frequent patterns list.

Product Barcode	Frequency
39	4321
48	3436
41	2024
32	1397
38	1348

5.3. Extracting frequent patterns list

After the construction of FUFP-Tree with pane-based partitioning sliding window approach, the association rules technique has been performed to discover the frequent patterns list. There are more than 8000 frequent patterns are extracted. The top 5 frequent itemsets are listed in Table 3.

5.4. Results analysis

The results of the proposed algorithm to construct the FUFP-Tree are exhibited. The results demonstrate that FUFP-Tree is very effective in terms of memory storage when finding correct frequent patterns from a streaming data. The runtime changes based on sliding windows configurations. As shown in Table 2, if there are more sliding windows, high number of incoming streaming data are processed in the same time, the runtime is reduced and the FUFP-Tree is built in short time.

[15] David del Rio Astorga, Manuel F. Dolz, Javier Fernández and J. Daniel García, "A generic parallel pattern interface for stream and data

6. Conclusion And Future Work

In this paper, the characteristics of streaming data have been discussed and presented some related algorithms for generating association rules from a dataset. The proposed technique works on speeding up streaming data mining using a sliding window technique to build FUFP-Tree with fault tolerance for any failure in the system. In the future, the proposed technique will be improved to adjust the number of needed sliding window and the pane size for each window dynamically according to the rate of streaming data. Such improvement results in memory usage reduction and speed up building FUFP-Tree.

REFERENCES

[1] A. Moustafa, Badr Abuelnasr, and Mohamed Said Abougabal, "Efficient mining fuzzy association rules from ubiquitous data streams", published in Alexandria Engineering Journal, vol. 54, pp. 163-174, June. 2015.

[2] B. Nath, D K Bhattacharyya, and A Ghosh, "Incremental Association Rule Mining: A Survey", published online in Wiley InterScience, vol. 3, pp. 157-169, February 2013.

[3] Luigi Troiano and Giacomo Scibelli, "Mining frequent itemsets in data streams within a time horizon", published online in ELSEVIER, vol. 89, pp. 21-37, January 2014.

[4] Akshita Bhandari, Ashutosh Gupta, and Debasis Das, "Improved apriori algorithm using frequent pattern tree for real time applications in data mining", published online in ELSEVIER, vol. 46, pp. 644-651, 2015.

[5] Jagrati Malviya, Anju Singh, and Divakar Singh, "An FP Tree based Approach for Extracting Frequent Pattern from Large Database by Applying Parallel and Partition Projection", published in International Journal of Computer Applications (0975 – 8887), vol. 114, pp. 1-5, March 2015.

[6] Hua-Fu Li, Man-Kwan Shan and Suh-Yin Lee, "DSM-FI: an efficient algorithm for mining frequent itemsets in data streams", published in Springer Knowledge and Information Systems, vol. 17, pp. 79-97, October 2008.

[7] Fabio Fumarola, Anna Ciampi, Annalisa Appice, and Donato Malerba "A Sliding Window Algorithm for Relational Frequent Patterns Mining from Data Streams", published in Springer International Conference on Discovery Science, vol. 5808, pp. 385-392, 2009.

[8] R. Agrawal, S. Ghosh, T. Imielinski, B. Iyer and A. Swami, "An Interval Classifier for Database Mining Applications," Proceedings of the VLDB Conference, pp.560-573, 1992.

[9] J. Li et al. Semantics and Evaluation Techniques for Window Aggregates in Data Streams. In ACM SIGMOD Conference, Baltimore, MD, USA, June 2005.

[10] T. Brijis, Retail market basket data set, Workshop on Frequent Itemset Mining Implementations (FIMI'03), 2003.

[11] Philippe Fournier-Viger, Espérance Mwamikazi, Ted Gueniche1 and Usef Faghihi, "MEIT: Memory Efficient Itemset Tree for Targeted Association Rule Mining", published in Springer.

[12] International Conference on Advanced Data Mining and Applications, vol. 8347, pp. 95-106, 2013.

[13] Chun-Wei LIN, Tzung-Pei HONG and Wen-Hsiang LU, "Using the Structure of Prelarge Trees to Incrementally Mine Frequent Itemsets", published in Springer New Generation Computing, vol. 28, pp 5–20, January 2010.

[14] Nataliya Shakhovska, Roman Kaminsky, Eugen Zasoba and Mykola Tsiutsiura, "Association rules mining in big data", published in International Journal of Computing, vol. 17, pp. 25-32, 2018.

[15] David del Rio Astorga, Manuel F. Dolz, Javier Fernández and J. Daniel García, "A generic parallel pattern interface for stream and data processing", published online in Wiley InterScience, vol. 29, pp. 1-12, May 2017.

- [16] Rashed Salem, Jérôme Darmont and Omar Boussaïd, “Efficient incremental breadth-depth XML event mining”, published in 15th International Database Engineering & Applications Symposium, pp. 197-203, September 2011.
- [17] Xiuli Yuan, “An improved Apriori algorithm for mining association rules”, published in AIP Conference Proceedings, vol. 1820, pp. 080005-1–080005-6, March 2017.
- [18] S.Alagukumar, R.Lawrance, “A selective analysis of microarray data using association rule mining”, published online in ELSEVIER, vol. 47, pp. 3-12, 2015.
- [19] Min Chen, Shiwen Mao and Yunhao Liu, “Big Data: A Survey”, published in Springer Mobile Networks and Applications, vol. 19, pp. 171-209, April 2014.