

Sequential Pattern Mining from Multidimensional Sequence Data in Parallel

Mahdi Esmaili and Mansour Tarafdar

Abstract—Sequential pattern mining is applicable in a wide range of applications since many types of data sets are in a time related format. Besides mining sequential patterns in a single dimension, mining multidimensional sequential patterns can give us more informative and useful patterns. Due to the huge increase in data volume and also quite large search space, efficient solutions for finding patterns in multidimensional sequence data are nowadays very important. For this reason, developing a parallel algorithm is necessary. In this paper, we present a multidimensional sequence model and a parallel algorithm follows the level-wise approach and all participating processors or workers generate candidate sequence and count their supports independently. Simulation experiments show good load balancing and scalable and acceptable speedup over different processors and problem sizes.

Index Terms—Data Mining, Sequential Patterns, Sequence Data, Parallel Algorithms.

I. INTRODUCTION

Data mining has been defined as “The nontrivial extraction of implicit, previously unknown, and potentially useful information from data.” Mining frequent patterns or itemsets is a fundamental and essential problem in many data mining applications. These applications include the discovery of association rules, strong rules, correlations, sequential rules, episodes, multidimensional patterns, and many other important discovery tasks [1].

One of the most important data mining problems is discovery of frequently occurring patterns in sequence data [2]. There are many domains where sequence mining has been applied, which include analysis of telecommunication systems, discovering customer buying patterns in retail stores, analysis of web access databases, and mining DNA sequences and gene structures.

Usually, sequence patterns are associated with different circumstances, and such circumstances form a multiple dimensional space. For example, bank customer sequences are associated with credit, branch, age, and others. It is interesting and useful to mine sequential patterns associated with multidimensional information [3].

Most of the proposed pattern mining algorithms are a variant of Apriori. Apriori-based algorithms show good performance with sparse datasets such as market-basket data, where the frequent patterns are very short. However, with

dense data sets such as telecommunications, where there are many and long frequent patterns, performance of these algorithms degrades incredibly.

On the one side, parallel and distributed computing is expected to relieve current mining methods from the sequential bottleneck, providing the ability to scale to massive datasets, and improving the response time. Achieving good performance on today's multiprocessor systems is a non-trivial task. The main challenges include synchronization and communication minimization, work-load balancing, find good data layout and data decomposition, and disk I/O minimization, which is especially important for data mining [4].

On the other side, the basic single dimension can not satisfy the requirement of multi-attribute analysis, which is often the case in actual system practice. To address this problem, multidimensional sequence pattern mining is developed.

The most time consuming operation in the discovery process of sequential patterns is the computation of the frequency of the occurrences of interesting subsequences in the sequence database. However, the number of sequential patterns grows exponentially and the task of finding all sequential patterns requires a lot of computational resources, which make it an ideal candidate for parallel processing. However, sequential pattern mining suffers from a scalability problem in both memory use and computational time when a dataset size is large. To perform mining on large datasets, we propose to parallelize this task. We demonstrate that our parallel algorithm can effectively alleviate the scalability problem.

In this paper, we present a model for multidimensional sequence pattern and then propose a parallel algorithm for mining sequential patterns from multidimensional sequence data. The algorithm is based on the data parallelism strategy of pattern growth. The experimental results show that our parallel algorithms usually achieve good load balancing and scalability.

The rest of the paper is organized as follows: We describe the sequence mining and related work in Section 2. Section 3 describes the problem definition. Algorithm and an experimental study are presented in Section 4, 5, and we conclude in Section 6.

II. SEQUENCE MINING AND RELATED WORK

The problem of mining frequent patterns in a set of data sequences together with a few mining algorithms was first introduced in [2]. They also presented three algorithms for solving this problem.

Mahdi Esmaili is with the Department of Computer science, Islamic Azad University (Kashan branch), Kashan, Iran (phone: +98-361-5550055; fax: +98-361-5550056; e-mail: M.Esmaili@iaukashan.ac.ir).

Mansour Tarafdar is MS student in Islamic Azad University (Qazvin branch), Qazvin, Iran (e-mail: tarafdar.mansour@gmail.com).

The sequence mining task is to discover a sequence of attributes, shared across time among a large number of objects in a given database. Many studies have been contributed to the efficient mining of sequential patterns in the literature, most of which was focused on developing efficient algorithms for finding all sequential patterns such as AprioriAll [2], GSP [5], SPADE [6], PrefixSpan [7] and so on.

In addition, enormous sizes of available databases and possibly large number of mined sequential patterns demand efficient and scalable parallel algorithms. While parallel association mining has attracted wide attention there has been relatively less work on parallel mining of sequential patterns [8]. Several parallel algorithms such as HPSPM [9], pSPADE [5], TPF [10] were proposed and have good performance. Moreover, recent research on sequential pattern mining has progressed closed sequential pattern mining, which can greatly reduce the number of frequent subsequences and improve the efficiency. Yet, the above work all assumes that the database is static, and a database updates requires rediscovering all the patterns by scanning the entire old and new database. Then, there is a need for efficient algorithms to update, maintain and manage the information discovered. Some incremental mining algorithms of sequential patterns were proposed and perform significantly better than the original approach of mining the whole updated database from scratch.

Other work contributes on an extension of the problem of sequential pattern mining like constraint-based in sequential pattern mining, mining cyclically repeated patterns, approximate mining of consensus sequential patterns, mining top-k closed sequential patterns, mining frequent max sequential patterns, mining long sequential patterns in a noisy environment, mining hybrid sequential patterns and sequential rules, etc.

Mining multidimensional sequential pattern is also one of the central topics in sequence mining as showed by the research efforts produced in recent years [3][11][12][13].

III. PROBLEM DEFINITION

There is a rich variety of sequence terminology but this section defines the couple of concepts according to our model. A table T is a set of tuples $\langle D, A_1, A_2, \dots, A_n \rangle$, where D is an attribute whose its domain is totally ordered. A sequence s is denoted by an ordered list $\langle t_1, t_2, \dots, t_k \rangle$, where t_i is a tuple, i.e., $D(t_1) \leq D(t_2) \leq \dots \leq D(t_k)$ for $1 \leq i \leq k$ and $D(t_i) = \text{value of } D \text{ tuple } t_i$. Every tuple has n analysis attributes along with an ordered value. A set of analysis attribute values can occur at most once in a same value of D , but can occur multiple times in different values of D attribute. The number of values of D in a sequence is called the length of the sequence. If the length of s is k , then we call it a k -sequence.

A sequence $S_1 = \langle a_1 a_2 \dots a_n \rangle$ is called a subsequence of another sequence $S_2 = \langle b_1 b_2 \dots b_m \rangle$ and S_2 is a super sequence S_1 , if there exist integers $1 \leq j_1 \leq j_2 \leq \dots \leq j_n \leq m$ such that $a_1 \subseteq b_{j_1}$, $a_2 \subseteq b_{j_2}, \dots, a_n \subseteq b_{j_n}$.

A multidimensional sequence database is of schema $\langle D, A_1, \dots, A_n, R_1, \dots, R_m \rangle$, where R_i are called relevant dimensions.

The schema is partitioned according to relevant attribute values and support computed by number of partitions that contain sequence. As mentioned previously, every partition is similar to table T , a set of tuples $\langle D, A_1, \dots, A_n \rangle$.

Given a minimum support threshold min-support , a multidimensional sequence S is called a multidimensional sequential pattern if and only if $\text{support}(S) \geq \text{min-support}$.

The problem of sequential pattern mining is to find the complete set of frequent sequential patterns satisfying a minimum support in the sequence database.

In this framework, $\langle (t_1, t_2), (t_3), (t_4, t_5, t_1) \rangle$ is a multidimensional sequence. All tuples in a bracket are assumed to occur at the same time. Due to t_1 and t_2 have the same value of D ; both of them appear in one bracket. A sequence is an ordered list of brackets in which ordered according to their associated time in a sequence. Number of brackets show length of sequence. Therefore, above sequence is a sequence with length 3. It is important to note that, each tuple can be repeated in different bracket (time), such as tuple t_1 , which appear in first and third bracket. On the other side, tuple ordering in every bracket is not important.

IV. THE ALGORITHM

There are two main paradigms that may be utilized in the implementation of parallel sequence mining: a data parallel approach or a task parallel approach. In data parallelism P processors work on distinct portions of the database, but synchronously process the global computation tree. In task parallelism, the processors share the database, but work on different classes in parallel.

In this section, we present a model for multidimensional sequence and describe a parallel algorithm based on data parallelism. This model originally introduced in [12]. We assume that we want to extract all multidimensional sequence patterns from a relational table T . We also assume that T contains at least one dimension whose domain is totally ordered, corresponding to the time dimension. Other dimensions divided to tree groups: analysis, relevant and, irrelevant dimensions. Irrelevant dimensions, as the name suggests, do not any role in my approach and eliminated from table T in preprocessing phase. Tuples over analysis dimensions are those that appear in the sequential patterns to be mined. The table is partitioned according to tuple values over relevant dimensions and support computed by number of partitions that contain sequence. Note that the length of a partition is defined as the number of distinct values of time dimensions. In almost all the test cases, we found that using this length for distributing partitions among processors make a good load balancing.

The well-known Apriori property holds on the algorithm, which states that any superpattern of an infrequent pattern is infrequent and any subpattern of a frequent pattern is also frequent. A dynamic programming approach can be used to generate sequence with length n from sequence with length $n-1$.

Fig.1 shows that how we can distribute 12 subsets among 4 processors based on sequence length. As you see, subsets ordered according to its length and allocate them to

processors. For example, subsets with length 6, 4, and 3 allocate to the first processor and subsets with length 5, 5, and 2 allocate to the fourth processor.

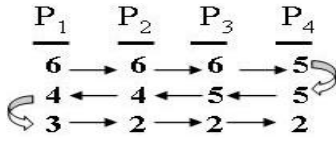


Fig. 1 Distribution among Processors

Given a table with length m and d_i is number of tuples with time dimension value i . Number of sequences with length 1, 2, 3, m computed as follows:

$$1 - \text{Sequence} : \sum_{i=1}^m (2^{d_i} - 1) \quad (1)$$

$$2 - \text{Sequence} : \sum_{i=1}^{m-1} (2^{d_i} - 1) \sum_{j=i+1}^m (2^{d_j} - 1) \quad (2)$$

$$3 - \text{Sequence} : \sum_{i=1}^{m-2} (2^{d_i} - 1) \sum_{j=i+1}^{m-1} (2^{d_j} - 1) \sum_{k=j+1}^m (2^{d_k} - 1) \quad (3)$$

$$m - \text{Sequence} : \prod_{i=1}^m (2^{d_i} - 1) \quad (4)$$

As you observe, search space is quite large and the serial algorithms are not scalable. To overcome this problem, it is necessary to implement scalable parallel algorithms of sequence mining algorithms.

The main steps of our parallel algorithm show in Fig. 2. Our proposed algorithm conducts a level-by-level candidate generation and test pruning following the Apriori property in workers independently. At each level, only potentially frequent candidates are generated and tested.

Fig. 2 Main Steps of Algorithm

V. EXPERIMENTS

1. One worker (processor) known as coordinator, partitions table T into P subsets according to relevant attributes.
2. Coordinator distributes P subsets among workers based on subset's length.
3. Every worker mines all k -sequence and computes its support (local support) and then sends to coordinator
4. Coordinator collects results and discards some sequences in which its support is less than min-support and sends sequence patterns to workers.
5. Repeat 3 and 4 until there is no sequence to be mined.

This section describes a set of simulations used to demonstrate the benefit of our approach. We simulate many experiments with various synthetic datasets on single processor machine. For this reason, we did not obtain excellent speedup but results show good scalability of our model, run time grows almost linearly when the database size increases.

Maximum length of patterns is 10 and even with corresponding value of length each examination took long time in this situation. Minimum value of support, number of analysis attributes, and size of database are variables which we have changed and obtain some plots. Note that we are looking for all of multidimensional sequential patterns.

Fig. 3 and 4 show data characteristics. Former shows number of patterns based on support value and later shows length and number of patterns in generated datasets.

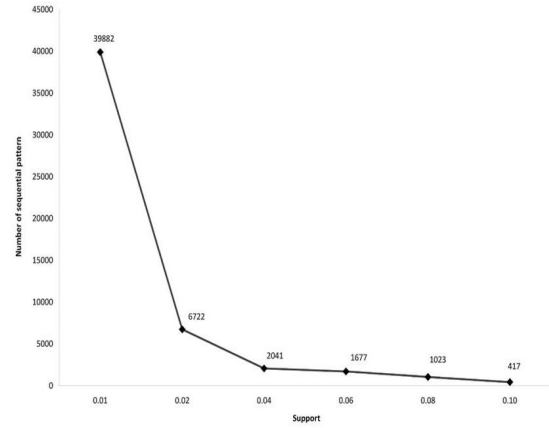


Fig. 3 Number of Patterns over Support

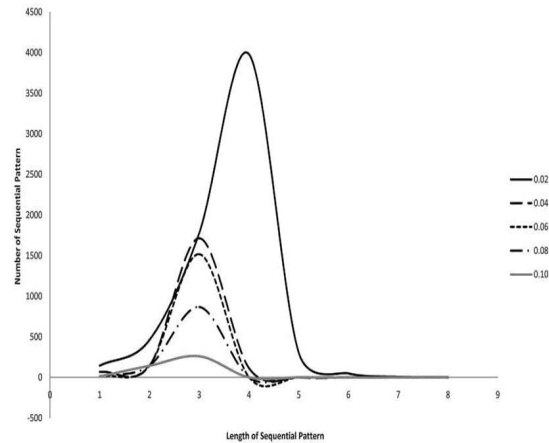


Fig. 4 Distribution of Datasets

Fig. 5 and 6 show the effect of the four processors over support value and database increases, respectively. These Figures show experimental performance results we obtained by simulating our model using up to 4 workers.

We can see how the system behavior is scalable. Speedup on processors is about 3.6 and run time decreases from 14605 to 4058 seconds.

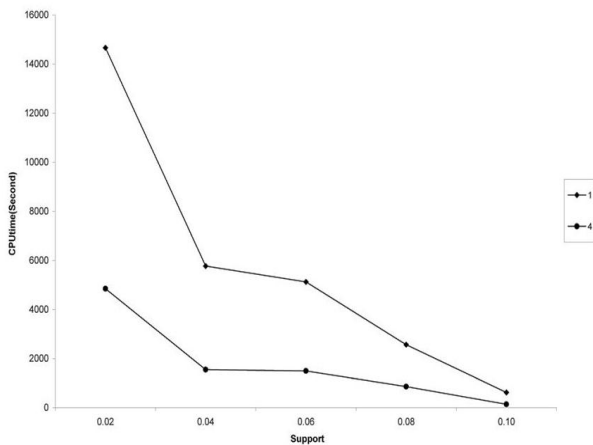


Fig. 5 Runtime over Support

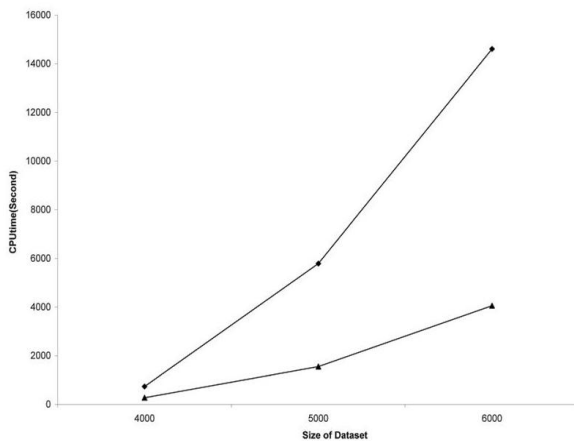


Fig. 6 Runtime over Database Size

VI. CONCLUSION

Multidimensional mining has been attracting attention in recent research into data mining. Very large search space and data volume have made many problems for serial algorithms to mine sequential patterns. In order to effectively mine, efficient parallel algorithm is necessary. In this paper, we theoretically present a multidimensional sequence model and then use SPMD (Single Program Multiple Data) strategy to parallelize multidimensional sequential pattern mining. According to this method a set of processors execute in parallel the same algorithm on different partitions of a dataset. The main goal of the algorithm is balanced workload among the processors and good scalability. We have simulated our parallel algorithm using several data sets on single processor. Simulation results on synthetic data sets show that if we implement this model on real computing environment, we may achieve good speedup.

The future research issues in multi-dimensional sequential pattern mining can be: multi-dimensional sequential pattern mining with constraints, interactive multi-dimensional sequential pattern mining or multi-dimensional sequential pattern mining integrated with taxonomies and hierarchies.

ACKNOWLEDGMENT

This work was supported in part by Islamic Azad University (Kashan branch) and M. Esmaili Author would likes to express his appreciation to his sponsor.

REFERENCE

- [1] J. Han, and M. Kamber, *Data Mining: Concepts and Techniques*, Academic Press, New York, 2001.
- [2] R. Agrawal, R. Srikant, *Mining Sequential Patterns*, Proc. of the 11th Int'l Conference on Data Engineering, 1995.
- [3] J. Han, J. Pei, K. Wang, Q. Chen, and U. Dayal, *Multi-dimensional sequential pattern mining*, Conference on Information and Knowledge Management, pp: 81-88, 2001 .
- [4] M. J. Zaki, C.-T. Ho, *Large scale parallel data mining*, Lecture notes in computer science, Vol 1759: Lecture notes in artificial intelligence, Springer, 2000 .
- [5] R. Srikant, R. Agrawal, *Mining sequential patterns: generalizations and performance improvements*, In Proc. of the 5th Int'l Conference on Extending Database Technology, pp: 3-17, 1996 .
- [6] M. J. Zaki, *An efficient algorithm for mining frequent sequences*, Machine Learning Journal, vol. 42(1/2), pp: 31-60, 2001 .
- [7] J. Pei, J. Han, H. Pinto, Q. Chen, U. Dayal, MC. Hsu, *PrefixSpan: Mining sequential patterns efficiently by prefix-projected pattern growth*, In Proc. of Int'l Conference on Data Engineering, pp: 215-224, 2001.
- [8] W. Jinlin, X. Chen, Z. Kefa, and W. Wei, *Parallel Research of Sequential Pattern Data Mining Algorithm*, Int'l Conference on Computer Science and Software Engineering, 2008 .
- [9] T. Shintani, M. Kitsuregawa, *Mining algorithms for sequential patterns in parallel: Hash based approach*, In Proc of the Second Pacific-Asia Conf on Knowledge Discovery and Data mining, pp: 283-294, 1998 .
- [10] V. Guralnik, N. Garg, G. Karypis, *Parallel tree projection algorithm for sequence mining*, Proc of 7th European Conf. on Parallel Computing, pp:310-320, Springer, 2001 .
- [11] S. de Amo, D. A. Furtado, A. Giacometti, and D. Laurent, *An apriori-based approach for first-order temporal pattern mining*, In Simposio Brasileiro de Bancos de Dados, 2004.
- [12] Plantevit M., Choong Y.W., Laurent A., Laurent D., Teisseire M. , *M2SP: Mining Sequential Patterns Among Several Dimensions*, Principles of Knowledge Discovery in Databases, PKDD, Volume 3721, pp: 205-216, 2005 .
- [13] C.-C. Yu, and Y.-L. Chen, *Mining Sequential Patterns from Multidimensional Sequence Data*, IEEE Transactions on Knowledge and Data Engineering archive, Volume 17, Issue 1, pp: 136-140, 2005.

Mahdi Esmaili is currently PhD student at University of Debrecen in Hungary. He received a Bs in software engineering from the University of Isfahan, and his MS in computer science at 1999 in Iran. He is a faculty member of Islamic Azad University (Kashan branch). He has taught in the areas of database and data mining and his research interests include data base, data mining and knowledge discovery, and recent research focusing on the end use of data mining. He has published books about database and information retrieval and also has published articles in international conferences.