Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

A Novel Approach for Scalability -Two Way Sequential Pattern Mining using UDDAG Dr.P.Raguraman¹ ¹ SRI SANKARA ARTS AND SCIENCE COLLEGE *Received: 9 December 2012 Accepted: 3 January 2013 Published: 15 January 2013*

7 Abstract

Traditional pattern growth-based approaches for sequential pattern mining derive length- (k +8 1) patterns based on the projected databases of length-k patterns recursively. At each level of 9 recursion, they unidirectionally grow the length of detected patterns by one along the suffix of 10 detected patterns, which needs k levels of recursion to find a length-k pattern. In this paper, a 11 novel data structure, UpDown Directed Acyclic Graph (UDDAG), is invented for efficient 12 sequential pattern mining. UDDAG allows bidirectional pattern growth along both ends of 13 detected patterns. Thus, a length-k pattern can be detected in $|\log 2 k + 1|$ levels of recursion 14 at best, which results in fewer levels of recursion and faster pattern growth. When minSup is 15 large such that the average pattern length is close to 1, UDDAG and PrefixSpan have similar 16 performance because the problem degrades into frequent item counting problem. However, 17 UDDAG scales up much better. It often outperforms PrefixSpan by almost one order of 18 magnitude in scalability tests. UDDAG is also considerably faster than Spade and 19 LapinSpam. Except for extreme cases, UDDAG uses comparable memory to that of 20 PrefixSpan and less memory than Spade and LapinSpam. Additionally, the special feature of 21 UDDAG enables its extension toward applications involving searching in large spaces. 22

23

Index terms— data mining algorithm, directed acyclic graph, performance analysis, sequential pattern, transaction database.

²⁶ 1 Introduction

EQUENTIAL pattern mining is an important data mining problem, which detects frequent subsequences in a 27 sequence database. A major technique for sequential pattern mining is pattern growth. Traditional pattern 28 growth-based approaches (e.g., PrefixSpan) derive length-(k + 1) patterns based on the projected databases 29 of a length-k pattern recursively. At each level of recursion, the length of detected patterns is grown by 1, and 30 patterns are grown unidirectionally along the suffix direction. Consequently, we need k levels of recursion to mine 31 a length-k pattern, which is expensive due to the large number of recursive database projections. In this paper, a 32 new approach based on UpDown Directed Acyclic Graph (UDDAG) is proposed for fast pattern growth. UDDAG 33 34 is a novel data structure, which supports bidirectional pattern growth from both ends of detected patterns. With 35 UDDAG, at level i recursion, we may grow the length of patterns by 2i_1 at most. Thus, a length-k pattern 36 can be detected in $|\log 2 |k + 1|$ levels of recursion at minimum, which results in better scale-up property for 37 UDDAG compared to PrefixSpan. Our extensive experiments clearly demonstrated the strength of UDDAG with its bidirectional pattern growth strategy. When minSup is very large such that the average length of patterns 38 is very small (close to 1), UDDAG and PrefixSpan have similar performance because in this case, the problem 39 degrades into a basic frequent item counting problem. However, UDDAG scales up much better compared to 40 PrefixSpan. It often outperforms PrefixSpan by one order of magnitude in our scalability tests. UDDAG is also 41 considerably faster than two other representative algorithms, Spade and LapinSpam. Except for some extreme 42

43 cases, the memory usage of UDDAG is comparable to that of PrefixSpan. UDDAG generally uses less memory

than Spade and LapinSpam. UDDAG may be extended to other areas where efficient searching in large searching spaces is necessary.

46 **2** II.

47 **3** Related Work

The problem of sequential pattern mining was introduced by Agrawal and Srikant [1]. Among the many algorithms 48 proposed to solve the problem, GSP ??17] and PrefixSpan[13], [14] represent two major types of approaches: a 49 prioribased and pattern growth-based. A priori principle states that any supersequence of a nonfrequent sequence 50 51 must not be frequent. A priori based approaches can be considered as breadth-first traversal algorithms because 52 they construct all length-k patterns before constructing length-(k+1) patterns. The AprioriAll algorithm [1] is 53 one of the earliest a prioribased approaches. It first finds all frequent item sets, transforms the database so that 54 each transaction is replaced by all frequent item sets it contains, and then finds patterns. The GSP algorithm [16] is an improvement over AprioriAll. To reduce candidates, GSP only creates a new length-k candidate when there 55 are two frequent length-(k _1) sequences with the prefix of one equal to the suffix of the other. To test whether 56 a candidate is a frequent length-k pattern, the support of each length-k candidate is counted by examining all 57 the sequences. The PSP algorithm ??12] is similar to GSP except that the placement of candidates is improved 58 through a prefix tree arrangement to speed up pattern(D D D D D D D D) 59 discovery. The SPIRIT algorithm [9] uses regular expressions as constraints and developed a family of 60

algorithms for pattern mining under constraints based on a priori rule. The SPaRSe algorithm [3] improves
 GSP by using both candidate generation and projected databases to achieve higher efficiency for high pattern
 density conditions.

64 **4** III.

Problem Definition a) Updown Directed Acyclic Graph-Based Sequential Pattern Mining UDDAG-based pattern mining approach, which first transforms a database based on frequent item sets, then partitions the problem, and finally, detects each subset using UDDAG. The absolute support for an item set in a sequence database is the number of tuples whose sequences contain the item set. An item set with a support larger than minSup is called a frequent item (FI) set. Based on frequent item sets, we transform each sequence in a database D into an alternative representation.

⁷¹ 5 i. Transformed Database Definition

⁷² Let D be a database and P be the complete set of sequential patterns in D, D' be its transformed database, ⁷³ substituting the ids of each item pattern contained in D' with the corresponding item sets, and denoting the ⁷⁴ resulted pattern set by P', we have P = P'.

⁷⁵ Based on frequent item sets, we transform each sequence in a database D into an alternative representation.

Steps involved in Database Transformation: 1. Find the set of frequent items in D. 2. Assign a unique id to
each FI in D and then replace each item set in each sequence with the ids of all the FIs contained in the item
For the database in Table 1, the FIs are: (1),(??), (3), (??), (??), (??), (1,2), (2,3).

By assigning a unique id to each FI, e.g., (1)-1, (1,2)-2, (2)-3, (2,3)-4, (3)-5, (4)-6, (5)-7, (6)-8, we can transform the database as shown in Table 2 (infrequent items are eliminated).

ii. Problem Partitioning Definition Let $\{x1, x2, \ldots, xt\}$ be the frequent item sets in a database D, $x1 < x2 < \ldots < xt$, the complete set of patterns (P) in D can be divided into t disjoint subsets. The ith subset (denoted by P xi, 1 < = i <= t) is the set of patterns that contains xi and FIs smaller than xi.

⁸⁴ iii. (Projected database) collection of all the tuples whose sequences contain an item set in a database D is ⁸⁵ called x-projected denoted by x The total number of different Frequent Items FIs in the Sequential Database are ⁸⁶ found by Database Transformation module. For a database with n different frequent items, its patterns can be ⁸⁷ divided into n disjoint subsets. The subset (1 < i < n) is the set of patterns that contain (the root item of the ⁸⁸ subset) and items smaller than i.Each subset i of the problem is mapped into a projected database denoted by (

i D).
P is partitioned into eight subsets as there are 8 FIs in table-2, the one contains 1 (P1), the one contains 2
and smaller ids (P2), . . . ,and the one contains 8 and smaller ids (P8).

Given the following database (P8) alone is found as given by 8

93 6 iii. UDDAG based Pattern Mining

In the ith subset, each pattern in the projected database (x D) can be divided into two parts, prefix and suffix of i.

The collection of all the prefix/suffix tuples of a frequent item set X in x D is called the prefix/suffixprojected database of x, denoted by Pre(x D) / Suf(x D).

To detect the sequential pattern in projected database (x D) P x, Detect patterns in Pre(x D) called pattern prefix (PP). Detect pattern in Suf(x D) called pattern suffix (PS) The above steps are repeated recursively until no frequent items are found in the pre(x D) / suf(x D). Combine the patterns of all the to derive P x.

The complete set of patterns the union of patterns of the all subsets or projected database (x D) detected above.

103 The Apriori property is used to reduce the number of candidate sets to be considered

¹⁰⁴ 7 Example for Pattern Mining

¹⁰⁵ 8 Performance Evaluation

We conducted an extensive set of experiments to compare our approach with other representative algorithms. All 106 the experiments were performed on a windows Server 2003 with 3.0 GHz Quad Core Intel Xeon Server and 16 107 GB memory. The algorithms we compared are PrefixSpan, Spade, and LapinSpam, which were all implemented 108 in C++ by their authors (Minor changes have been made to adapt Spade to Windows). Two versions of UDDAG 109 were tested. UDDAG-by uses bit vector to verify candidates and UDDAG-co uses co-occurrences to verify 110 candidates whenever possible. We perform two studies using the same data generator as in [14]: 1) Comparative 111 study, which uses similar data sets as that in [14]; 2) Scalability study. The data sets were generated by 112 maintaining all except one of the parameters shown in Table ?? fixed, and exploring different values for the 113 remaining ones. 114

115 V.

116 9 Conclusion

In this paper, a novel data structure UDDAG is invented for efficient pattern mining. The new approach grows patterns from both ends (prefixes and suffixes) of detected patterns, which results in faster pattern growth because of less levels of database projection compared to traditional approaches. Extensive experiments on both comparative and scalability studies have been performed to evaluate the proposed algorithm.

One major feature of UDDAG is that it supports efficient pruning of invalid candidates. This represents a promising approach for applications involving searching in large spaces. Thus, it has great potential to related areas of data mining and artificial intelligence. In the future, we expect to further improve UDDAG-based pattern mining algorithm as follows: 1) Currently, FI detection is independent from pattern mining. Practically, the knowledge gained from FI detection may be useful for pattern mining. In the future, we will integrate the solutions of the two so that they can benefit from each other. 2) Different candidate verification strategies may have different impacts to the efficiency of the algorithm. In the future, we will study more efficient verification

128 strategy. 3) UDDAG has big impact to the memory usage when the number of patterns in a subset is extremely 129 large. In the future, we will find an efficient way to store UDDAG.

 $^{^1 \}ensuremath{\mathbb C}$ 2013 Global Journals Inc. (US) Global Journal of Computer Science and Technology



RESEARCH | DIVERSITY | ETHICS

Figure 1: C 2 .

Seq. Id	Sequence
1	<1 (1,2,3) (1,3) 4 (3,6)>
2	<(1,4) 3 (2,3) (1,5)>
3	<(5,6) (1,2) (4,6) 3 2>
4	<57 (1,6) 3 2 3>

Figure 2: Assuming 8

Seq. Id	Sequence
1	<1 (1,2,3,4,5) (1,5) 6 (5,8)>
2	<(1,6) 5 (3,4,5) (1,7)>
3	<(7,8) (1,2,3) (6,8) 5 3>
4	<7 (1,8) 5 3 5>

Figure 3: C

1

 $\mathbf{2}$

 $\mathbf{2}$

Figure 5: Table 2 :

9 CONCLUSION

- [Berkovich et al. ()] 'A Bit-Counting Algorithm Using the Frequency Division Principle'. S Berkovich , G Lapir
 , M Mack . Software: Practice and Experience 2000. 30 (14) p. .
- [Chen and Xiao] 'BISC: A Binary Itemset Support Counting Approach Towards Efficient Frequent Itemset
 Mining'. J Chen , K Xiao . ACM Trans. Knowledge Discovery in Data
- IGrahne and Zhu ()] 'Efficiently Using Prefix-Trees in Mining Frequent Itemsets'. G Grahne , J Zhu . Proc.
 Workshop Frequent Itemset Mining Implementations (FIMI '03), (Workshop Frequent Itemset Mining
 Implementations (FIMI '03)) 2003.
- [Agrawal and Srikant ()] 'Fast Algorithms for Mining Association Rules'. R Agrawal, R Srikant. Proc. 20th Int'l Conf. Very Large Data Bases (VLDB), (20th Int'l Conf. Very Large Data Bases (VLDB)) 1994. p. .
- [Antunes and Oliveira ()] 'Generalization of Pattern-Growth Methods for Sequential Pattern Mining with Gap
 Constraints'. C Antunes , A L Oliveira . Proc. Int'l Conf. Machine Learning and Data Mining, (Int'l Conf.
 Machine Learning and Data Mining) 2003. 2003. p. .
- [Chen and Cook (2007)] 'Mining Contiguous Sequential Patterns from Web Logs'. J Chen, T Cook. Proc. World
 Wide Web Conf. (WWW '07) Poster Session, (World Wide Web Conf. (WWW '07) Poster Session) May
 2007.
- [Agrawal and Srikant ()] 'Mining Sequential Patterns'. R Agrawal , R Srikant . CHEN:
 ANUPDOWNDIRECTEDACYCLICGRAPHAPPROACHFOREQUENTIALPATTERNMINING927 Proc. Int'l Conf.
 Data Eng. (ICDE '95), (Int'l Conf. Data Eng. (ICDE '95)) 1995. p. .
- 148 [Ayres et al. ()] 'Sequential Pattern Mining Using a Bitmap Representation'. J Ayres , J Gehrke , T Yu , J
- Flannick . Proc. Int'l Conf. Knowledge Discovery and Data Mining, (Int'l Conf. Knowledge Discovery and Data Mining) 2002. 2002. p. .
- 151 [Garofalakis et al. ()] 'SPIRIT: Sequential Pattern Mining with Regular Expression Constraints'. M Garofalakis
- 152 , R Rastogi , K Shim . Proc. Int'l Conf. Very Large Data Bases (VLDB '99), (Int'l Conf. Very Large Data
- 153 Bases (VLDB '99)) 1999. p. .