Global Journals LATEX JournalKaleidoscopeTM

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.*

OSSM: Ordered Sequence Set Mining for Maximal Length Frequent Sequences A Hybrid Bottom-Up-Down Approach

Anurag Choubey¹, Dr. Ravindra Patel² and Dr. J.L. Rana³

¹ Rajiv Gandhi Technological University, Bhopal, India

Received: 13 April 2012 Accepted: 2 May 2012 Published: 15 May 2012

7 Abstract

3

Δ

5

16

The process of finding sequential rule is indispensable in frequent sequence mining. Generally, 8 in sequence mining algorithms, suitable methodologies like a bottom-up approach is used for 9 creating large sequences from tiny patterns. This paper proposed on an algorithm that uses a 10 hybrid two-way (bottom-up and top-down) approach for mining maximal length sequences. 11 The model proposed is opting to bottom-up approach called "Concurrent Edge Prevision and 12 Rear Edge Pruning (CEGREP)" for itemset mining and top-down approach for maximal 13 length sequence mining. It also explains optimality of top-to-bottom approach in deriving 14 maximal length sequences first and lessens the scanning of the dataset. 15

17 Index terms—

18 1 Introduction

esearchers feel enthusiastic on the sequential pattern mining problems and wide range of possibilities of 19 applications regarding the envisaging of the customer buying patterns and scientific discoveries [1,2,3,4,5] 20 discussed by Agrawal and Srikant [1]. Let us explain with an example like finding the given time stamped 21 sequences of purchase made by a customer. In this example the main objective is to find sequence of same time 22 stamped list of items purchased by the customer. So the algorithm of sequence pattern mining should concentrate 23 on finding the repeated sequences which are called as frequent sequence. Such sequences list out the frequency 24 of common occurrences. several heuristics like GSP [1], SPADE [3], Prefix Span [2] and the SPIRIT [4] attempt 25 to find the frequent patterns in productive method by striving to cut short a series, hence decrease search space. 26 The GSP algorithm [1] utilizes the anti-monotone property (all subsequences of a frequent sequence are also 27 frequent). 28

The SPADE finds frequent sequences using the lattice search [3] and intersection based approach. In this 29 particular method the sequence database is converted into a vertical format. The candidate sequences will 30 be made into different groups. These frequent sequences will be listed in SPADE utilizing two methods namely 31 breadth first method and depth first method. The base is calculated for the produced sequences. The approach of 32 mentioned three algorithms can be grouped as the candidate-production with a base evaluation. The PrefixSpan 33 [2] algorithm adopts growth method pattern. It utilizes recorded database to accomplish. Prefix is Projected 34 Sequential Pattern mining which checks prefix subsequences and includes the postfix sub sequences into the 35 databases. 36

37 **2** II.

38 3 Related work

The sequential item set mining problem was initiated by Agrawal and Srikant, and the same developed a filtered algorithm, GSP [1], based on the Apriori property [1]. Since then, lots of sequential item set mining algorithms are being developed to increase the efficiency. Some are SPADE [3], PrefixSpan [2], and SPAM [11]. SPADE

42 [3] is on principle of vertical id-list format and it uses a lattice-theoretic method to decompose the search space

into many tiny spaces, on the other hand PrefixSpan [2] implements a horizontal format dataset representation 43 and mines the sequential item sets with the pattern-growth paradigm: grow a prefix item set to attain longer 44 sequential item sets on building and scanning its database. The SPADE [3] and the PrefixSPan [2] highly perform 45 GSP [1]. SPAM [11] is a recent algorithm used for mining lengthy sequential item sets and implements a vertical 46 bitmap representation. Its observations reveal, SPAM [11] is more efficient in mining long item sets compared 47 to SPADE [3] and PrefixSpan [2] but, it still takes more space than SPADE [3] and PrefixSpan [2]. Since the 48 frequent closed item set mining [12], many capable frequent closed item set mining algorithms are introduced, like 49 A-Close [12], CLOSET [13], CHARM [14], and CLOSET+ [15]. Many such algorithms are to maintain the ready 50 mined frequent closed item sets to attain item set closure checking. To decrease the memory usage and search 51 space for item set closure checking, two algorithms, TFP [17] and CLOSET+2, implement a compact 2-level 52 hash indexed result-tree structure to keep the readily mined frequent closed item set candidates. Some pruning 53 methods and item set closure verifying methods, initiated that can be extended for optimizing the mining of closed 54 sequential item sets also. CloSpan is a new algorithm used for mining frequent closed sequences [16]. It goes by 55 the candidate maintenance-and-test method: initially create a set of closed sequence candidates stored in a hash 56 indexed result-tree structure and do post-pruning on it. It requires some pruning techniques such as Common 57 58 Prefix and Backward Sub-Itemset pruning to prune the search space as CloSpan [16] requires maintaining the 59 set of closed sequence candidates, it consumes much memory leading to heavy search space for item set closure 60 checking when there are more frequent closed sequences. Because of which, it does not scale well the number of 61 frequent closed sequences.

62 **4 III.**

63 Sequence, sub sequence and frequent sequences

We can say a sequence means an ordered set of events [1], set of events S is said to subsequence of 2 S. The sequence database S is a set of the form (tid, s) where tid is the transaction-id and s is the sequence generated from transaction.

Let the minimum support as a threshold defined by user which indicates the desired minimum occurrences of the sequence to be claimed as frequent.

A sequences j is lengthiest if ? ? s ! {s $S|i = \{1...m\}}$ j i

. We explain the maximal support of an event e which consists of items m i i MSF s MS e = = ?

The Maximal Support Factor is the threshold used in our proposal to minimize the subsequence search in top to bottom approach.

73 We apply an ordered search on sequence database, hence the sequence database will be ordered in descending 74 manner by MSF of sequences. Then the search for super sequence of a given sequence is limited to the sequences 75 with greater or equal MSF of the given sequence.

The preprocessed dataset with the sorted transaction list has several properties that can be used to cut short the search space which are hypothesized below: Ossm overview T is a transaction represented by transaction id tid and s I is an itemset that belongs to set I of frequent itemsets. ? Build sequences by grouping items as events in a given transaction. Here we follow the top-to-bottom approach to build events. First we build events based on maximal length itemsets and continue the process in descending order of the itemset lengths. $|| \max() \{ | ($ $1 \ 0 \} 1 \ SM$ Occurrence s i if s s i or i c c i i = ? = = ? = ? Here | | SM is set of

? Measure the Maximal support of each event e of the given transaction tid T (refer section 3 for measuring
 maximal support MS for an event e of transaction tid T).

? Measure the Maximal Support Factor MSF of each sequence tid S (Refer section 3 for measuring Maximal
 Support Factor MSF of sequence tid S)

? Order the sequences in sequence dataset S in descending manner of their MSF. ? Build weighted acyclic directed graph from frequent itemsets of length two. o Elements of the itemset considered as vertices o Support of that itemset considered as edge weight ? Apply WFI algorithm to find critical path between any two items which represents the maximal sequence between two elements opted. ? Apply all four properties that are hypothesized in section 3 to discard, prune sequences or select maximal length sequences.

91

V

OSSM: Top-Down Ordered Sequence Set Mining for maximal length sequences a) Concurrent Edge Prevision
 and Rear Edge Pruning (CEG&REP) [7] in OSSM i.

94 5 Preprocess

Dataset preprocessing and itemsets Database initialization is performed by us as the first stage of proposal. As we find itemsets with single element, we in parallel prune it with the itemsets of single elements if the support of the selected itemsets is less than the required support.

98 ii.

⁹⁹ 6 Concurrent Edge Prevision

100 In this phase, we select all itemsets from given itemset database as input in parallel. Then we start projecting 101 edges from each selected itemset to all possible elements. The first iteration includes the pruning process in parallel, from the second iteration onwards this pruning is not required, which we claimed as an efficient process compared to other similar techniques like BIDE [8]. In first iteration, we project an itemset p s that spawned

104 from selected itemset i s from S DB and an element i e considered from 'I'. iii.

105 7 Rear Edge Pruning

If any of '()() ts p ts p f s f s ? that edge will be pruned and all disjoint graphs except p s will be considered as closed sequence and moves it into S DB and remove all disjoint graphs from memory. The termination of above process do not take place till the graph becomes empty, i.e. till the elements which are connected through transitive edges and projecting itemsets are available in the memory.

¹¹⁰ 8 b) Building Sequence Set from Transaction Dataset

111 Here in this section we explore the process of building sequence dataset.

TD is the given transaction dataset I is the set of closed frequent itemsets of length 1 to m . Here m is the maximal length of the itemset.

114 ? Initially set of closed frequent itemsets is ordered in descending manner by itemset length. ? For each 115 transaction tid T in the given transaction dataset TD o Build events based on the closed frequent itemsets of the 116 set I such that the event lengths will be decided in the descending order of the frequent itemset length.

o Initially events with length of m determined that is maximal length of the frequent itemsets. Then As a part of the OSSM, we order the sequence dataset in the descending manner of the Maximal support Factor (refer section 3 for details).

Let S be the determined sequence dataset from the given transaction dataset TD and set of frequent itemsets I. In this phase of OSSM, we explore the building of a weighted acyclic directed graph. Let 2 I be the set of frequent itemsets of length 2.

Let G be the graph initially with vertex count of zero | | 0 V = and edge count of zero | | 0 E =. Here V is

vertex set and E is edge set. In this phase we apply WFI algorithm [9,10]. In the first pass of algorithm we try

to identify and evaluate potential long and rich candidates. The rich sequences are the one whose constituent

¹²⁶ 2-sequences have high support. In the directed graph, the 2-sequence frequencies are represented by the edge ¹²⁷ weights; we can easily compute the path with the highest weights between all pairs of nodes. Here we use WFI

weights; we can easily compute the path with the highest weights between all pairs of nodes. Here we use WF1 algorithm [9,10] for the purpose finding critical path(a path with maximal vertex count) with maximal weights.

¹²⁹ 9 f) Sequence Evolution

Each critical path generated from the graph G will be considered as candidate sequence and stored in candidate sequence set css g) Verifying Frequency of Candidate Sequence $12^{1/2}$

 $^{^{1}}$ © 2012 Global Journals Inc. (US)

 $^{^{2}}$ © 2012 Global Journals Inc. (US)



Figure 1: ?

1			
sequence 1 S is said to be a subsequence of an	oth	er	
sequence S	2	if	and

e 2 3 < <

i e is event of 2 S . The 2 S is said to super sequence of 1 S and 1

$$123Sssssj = j 1 n is j + n$$

$$+ s i j$$

$$=$$

$$?$$
And every event i
$$i s = i m i = e i e ? m , here$$

$$i$$

$$i$$

[Note: { , , ,.... , ,..... | ! {1, 2,3, , 1,.... }} s is considered as an item set, which is a non-empty, unordered, finite set of items, which can be represented as 1 2 3 { , , ,....., | ! {1, 2,3, 4,...., }} e i for e m =]

Figure 2:

0	tid i T ? and	$\{ \mid s \text{ i I for } s ? \}$	$\{1, \} I =$
o Here tid			

Figure 3: ?

? For each itemset s d i ? of 2 I build an edge ed in graph G o Let consider item

add item s to vertex set V . Increment $\mid \mid$ V by 1. o Let consider item

and add item d to vertex setV . Increment | | V by
o by 1.
o Add support of s d i ? as weight to edge ed .
e) Finding Critical Paths between two Items as

Maximal Length Sequences

Figure 4:

sds i??	as source vertex, and
sdd i??	as des- tination vertex,

o Find	() MSF cs			
o For { s S MSF s MSF cs () ()} ? ?	е	ach	sequence	s (refer
<pre>section 3) ? If cs s ? then increment sup() sc by 1 o If sup() cs st ? (refer the hypothesis 2 in section then</pre>	3)			esis lii
? ? For each candidate sequence				' cs o
sequence set css such that ' cs cs ? ? If	' cs cs ? the	en consider		date , _{cs} quent
hypothesis 3 in section 3) and move				' cs fi to
fss o Else 2 For each candidate sequence				,
: For each candidate sequence				date
sequence set css such that ' cs cs ? ? If	cs cs ?	' then o	consider	' cs as quent
prune it from css (refer hypothesis 4 in section 3) h) Finding Maximal Length Sequences				I
in previous phase ? is true	Let iss be t	he frequent	sequence set generated	
o For each frequent sequence { 'fs fss MSF fs (') ?			?	' fs suo ()} M
and sts is true ? If o	fs fs ?	' then s	et sts as false	

[Note: Finally, maximal length frequent sequence set mlfss contains sequences that are not subsequence of any other frequent sequences. Global Journal of Computer Science and Technology Volume XII Issue VII Version I]

Figure 5: ?

132 VI.

133 .1 Conclusion

The proposed ordered sequence set mining (OSSM) approach is a scalable and optimal because of its hybrid 134 bottom-up-down approach. OSSM is supported by our earlier Concurrent Edge Prevision and Rear Edge Pruning 135 (CEG&REP) [7] for frequent closed itemset mining, which was proven as efficient in memory usage and scalable 136 on dense datasets. A novel mechanism of sequence dataset generation from transaction dataset is introduced 137 in this paper. The proposed OSSM is capable to generate the longest candidate sequence by weighted acyclic 138 directed graph construction and also efficient and scalable to find frequent sequence set and maximal length 139 frequent sequence set due to top-down approach. To compute the support for a candidate sequence, it uses the 140 maximal support factor of sequences. And the order approach that ordering the sequence dataset in descending 141 order of MSF ensures that whole data set is not scanned. Also, if the data set contains long regular sequences, it 142 is identified early enough and thus all the subsequences of this is also marked regular and need not be evaluated. 143 The longest possible sequence is build up by bottom up algorithms starting from 2sequence. 144 [Warshall ()] 'A theorem on boolean matrices'. S Warshall . Journal of the ACM 1962. 9. 145

- [Agrawal and Srikant ()] R Agrawal, R Srikant. Mining Sequential Patterns: Generalizations and Performance
 Improvements Proceedings of the 5th International Conference on Extending Database Technology: Advances
 in Database Technology p, 1996. p. .
- [Floyd ()] 'Algorithm 97: Shortest path'. R Floyd . Communications of the ACM 1962. 5.
- [Antunes and Oliveira ()] Antunes, A Oliveira. Generalization of Pattern-Growth Methods for Sequential
 Pattern Mining with Gap Constraints" in Int'l Conf Machine Learning and Data Mining, 2003. p. .
- 152 [Aloy et al. ()] 'Automated Structure-based Prediction of Functional Sites in Proteins: Applications to Assessing
- the Validity of Inheriting Protein Function From Homology in Genome Annotation and to Protein Docking'.
 P Aloy, E Querol, F X Aviles, M J E Sternberg. Journal of Molecular Biology 2002. p. 311.
- [Wang and Han ()] BIDE: Efficient Mining of Frequent Closed Sequences, Jianyong Wang , Jiawei Han . 2004.
 p. .
- [Zaki and Hsiao (2002)] 'CHARM: An efficient algorithm for closed itemset mining'. M Zaki , C Hsiao . SDM'02,
 (Arlington, VA) April 2002.
- [Choubey et al. ()] 'Concurrent Edge Prevision and Rear Edge Pruning Approach for Frequent Closed Itemset
 Mining"; (IJACSA)'. Anurag Choubey , Dr , Dr J L Patel , Rana . International Journal of Advanced
 Computer Science and Applications 2011. 2 (11) .
- [Pei et al. (2002)] 'Constraint-based sequential pattern mining in large databases'. J Pei , J Han , W Wang .
 CIKM'02, (McLean, VA) Nov. 2002.
- [Pasquier et al. (1999)] 'Discoving frequent closed itemsets for association rules'. N Pasquier , Y Bastide , R
 Taouil , L . *ICDT'99*, (Jerusalem, Israel) Jan. 1999.
- [Kohavi et al. ()] R Kohavi , C Brodley , B Frasca , L Mason , Z Zheng . KDD-cup 2000 organizers' report:
 Peeling the Onion. SIGKDD Explorations, 2000. 2.
- [Garofalakis et al. ()] 'Mining Sequential Patterns with Regular Expression Constraints'. M Garofalakis , R
 Rastogi , K Shim . *IEEE Transactions on Knowledge and Data Engineering* 2002. 14 (3) p. .
- [Han et al. (2002)] 'Mining Top-K Frequent Closed Patterns without Minimum Support'. J Han , J Wang , Y
 Lu , P Tzvetkov . *ICDM'02*, (Maebashi, Japan) Dec. 2002.
- [Lin ()] 'Pincer Search: A new algorithm for discovering the maximum frequent set'. D I Lin , Kedem , Z .
 International conference on Extending database technology, 1998.
- [Han ()] 'PrefixSpan: Mining Sequential Patterns Efficiently by Prefix-Projected Pattern Growth'. Pei J Han , J
 Int'l Conf Data Engineering, 2001. p. .
- [Zaki ()] 'SPADE: An Efficient Algorithm for Mining Frequent Sequences'. M Zaki . Machine Learning, 2001.
 Kluwer Academic Pulishers. 42 p. .
- [Zaki (2001)] SPADE: An Efficient Algorithm for Mining Frequent Sequences, Machine Learning, v.42 n.1-2, M
 J Zaki . January-February 2001. p. .