

Mining Closed Itemsets for Coherent Rules: An Inference Analysis Approach

Mr. Kalli Srinivasa Nageswara Prasad¹ and Prof. S.Ramakrishna²

¹ Sri Venkateswara University, Tirupati.

Received: 1 September 2011 Accepted: 24 September 2011 Published: 6 October 2011

Abstract

Past observations have shown that a frequent item set mining algorithm are alleged to mine the closed ones because the finish offers a compact and a whole progress set and higher potency. Anyhow, the most recent closed item set mining algorithms works with candidate maintenance combined with check paradigm that is dear in runtime likewise as area usage when support threshold is a smaller amount or the item sets gets long. Here, we show, PEPP with inference analysis that could be a capable approach used for mining closed sequences for coherent rules while not candidate. It implements a unique sequence closure checking format with inference analysis that based mostly on Sequence Graph protruding by an approach labeled Parallel Edge projection and pruning in brief will refer as PEPP. We describe a novel inference analysis approach to prune patterns that tends to derive coherent rules. A whole observation having sparse and dense real-life information sets proved that PEPP with inference analysis performs larger compared to older algorithms because it takes low memory and is quicker than any algorithms those cited in literature frequently.

Index terms— Data mining, Association Rule Mining, Closed itemset, Frequent Itemset, KDD, PEPP.

1 INTRODUCTION

Association rule mining, introduced in [28], is considered as one of the most important tasks in Knowledge Discovery in Databases [29]. Among sets of items in transaction databases, it aims at discovering implicative tendencies that can be valuable information for the decision-maker. An association rule is defined as the implication $X \Rightarrow Y$, described by two interestingness measures support and confidence, where X and Y are the sets of items and $X \cap Y = \emptyset$.

Apriori [28] is the first algorithm proposed in the association rule mining field and many other algorithms were derived from it. It is very well known that mining algorithms can discover a prohibitive amount of association rules; Starting from a database, it proposes to extract all association rules satisfying minimum thresholds of support and confidence. For instance, thousands of rules are extracted from a database of several dozens of attributes and several hundreds of transactions. Furthermore, as suggested by Silbershatz and Tuzilin [30], valuable information is often represented by those rare low support and unexpected association rules which are surprising to the user. So, the more we increase the support threshold, the more efficient the algorithms are and the more the discovered rules are obvious, and hence, the less they are interesting for the user. As a result, it is necessary to bring the support threshold low enough in order to extract valuable information. Unfortunately, the lower the support is, the larger the volume of rules becomes, making it intractable for a decision-maker to analyze the mining result. Experiments show that rules become almost impossible to use when the number of rules overpasses 100. Thus, it is crucial to help the decisionmaker with an efficient technique for reducing the number of rules.

To overcome this drawback, several methods were proposed in the literature. On the one hand, different algorithms were introduced to reduce the number of itemsets by generating closed [31], maximal [32] or optimal

itemsets [33], and several algorithms to reduce the number of rules, using non redundant rules [34], [35], or pruning techniques [36]. On the other hand, post processing methods can improve the selection of discovered rules. Different complementary post processing methods may be used, like pruning, summarizing, grouping, or visualization [37]. Pruning consists in removing uninteresting or redundant rules. In summarizing, concise sets of rules are generated. Groups of rules are produced in the grouping process, and the visualization improves the readability of a large number of rules by using adapted graphical representations.

However, most of the existing post processing methods are generally based on statistical information in the database. Since rule interestingness strongly depends on user knowledge and goals, these methods do not guarantee that interesting rules will be extracted. In this paper, we propose a novel framework to identify closed itemsets. Associations are discovered based on inference analysis. The principle of the approach considers that an association rule should only be reported when there is enough interest gain claimed during inference analysis in the data. To do this, we consider both presence and absence of items during the mining. An association such as beer ? nappies will only be reported if we can also find that there are fewer occurrences of \neg beer ? nappies nappies. This approach will ensure that when a rule such as beer ? nappies is reported, it indeed has the strongest interest in the data as comparison was made on both presence and absence of items during the mining process.

2 II.

3 RELATED WORK

The sequential item set mining problem was initiated by Agrawal and Srikant, and the same was developed as a filtered algorithm, GSP [2], basing on the Apriori property [19]. Since then, lots of sequential item set mining algorithms are being developed for efficiency. Some are, SPADE [4], Prefixspan [5], and SPAM [6]. SPADE 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 implements a horizontal format dataset representation and mines the sequential item sets with the pattern-growth paradigm: grow a prefix item set to attain longer sequential item sets on building and scanning its database. The SPADE and the Prefixspan highly perform GSP. SPAM is a recent algorithm used for mining lengthy sequential item sets and implements a vertical bitmap representation. Its observations reveal, SPAM is better efficient in mining long item sets compared to SPADE and Prefixspan but, it still takes more space than SPADE and Prefixspan. Since the frequent closed item set mining [15], many capable frequent closed item set mining algorithms are introduced, like A-Close [15], CLOSET [20], CHARM [16], and CLOSET+ [18]. Many such algorithms are to maintain the ready mined frequent closed item sets to attain item set closure checking. To decrease the memory usage and search space for item set closure checking, two algorithms, TFP [21] and CLOSET+2, implement a compact 2-level hash indexed result-tree structure to keep the readily mined frequent closed item set candidates. Some pruning methods and item set closure verifying methods, initiated that can be extended for optimizing the mining of closed sequential item sets also. CloSpan is a new algorithm used for mining frequent closed sequences [17]. It goes by the candidate maintenance-and-test method: initially create a set of closed sequence candidates stored in a hash indexed result-tree structure and do post-pruning on it. It requires some pruning techniques such as Common Prefix and Backward Sub-Item set pruning to prune the search space as CloSpan requires maintaining the set of closed sequence candidates, it consumes much memory leading to heavy search space for item set closure checking when there are more frequent closed sequences. Because of which, it does not scale well the number of frequent closed sequences. BIDE [26] is another closed pattern mining algorithm and ranked high in performance when compared to other algorithms discussed. BIDE projects the sequences after projection it prunes the patterns that are subsets of current patterns if and only if subset and superset contains same support required. But this model is opting to projection and pruning in sequential manner. This sequential approach sometimes turns to expensive when sequence length is considerably high. In our earlier literature [27] we discussed some other interesting works published in recent literature.

4 III.

5 DATASET ADOPTION AND FORMULATION

Item Sets 'I': A set of diverse elements by which the sequences generate.

6 ?

Represents a sequence 's' of items those belongs to set of distinct items 'I', 'm' is total ordered items and $P(e_i)$ is a transaction, where e_i usage is true for that transaction. S: represents set of sequences, 't' represents total number of sequences and its value is volatile and s_j : is a sequence that belongs to S.

Subsequence: is a sequence p s of sequence set 'S' is considered as subsequence of another sequence q s of Sequence Set 'S' if all items in sequence S p is belongs to s q as an ordered list. This can be formulated as If $1 (() n_{pi} q p q i s s s s = ? ? ? ?$ Then $1 l : n m pi qj i j s s = < ? ?$ where $p q s S$

7 CLOSED ITEMSET DISCOVERY a) PEPP: Parallel Edge Projection and Pruning Based

Sequence Graph protrude [28] i. Preprocess:

As a first stage of the proposal we perform dataset preprocessing and itemsets Database initialization. We find itemsets with single element, in parallel prunes itemsets with single element those contains total support less than required support.

ii. Forward Edge Projection:

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

8 b) Description of Inference Analysis

Set $I = \{i_1, i_2, \dots, i_m\}$ be the universe of items composed of m different attributes, $i_k (k=1,2,\dots,m)$ is item. Transaction database D is a collection of transaction T , A transaction $t = (tid, X)$ is a tuple where tid is a unique transaction ID and X is an itemset. The count of an itemset X in D , denoted by $count(X)$, is the number of transactions in D containing X . The support of an itemset X in D , denoted by $supp(X)$, is the proportion of transactions in D that contain X . The negative rule $X \rightarrow \neg Y$ holds in the transaction set D with confidence $conf(X \rightarrow \neg Y) = \frac{supp(X \wedge \neg Y)}{supp(X)}$.

In Transaction database, each transaction is a collection of items involved sequences. The issue of mining association rules is to get all association rules that its support and confidence is respectively greater than the minimum threshold given by the user. The issues of mining association rules can be divide into two sub-issues as follows:

Find frequent itemsets, Generate all itemsets that support is greater than the minimum support. Generate association rules from frequent itemsets. In logical analysis, the direct calculation of support is not convenient, To calculate the support and confidence of negative associations using the support and confidence

9 VI. COMPARATIVE STUDY

This segment focuses mainly on providing evidence on asserting the claimed assumptions that 1) The PEPP is similar to BIDE which is actually a sealed series mining algorithm that is competent enough to momentarily surpass results when evaluated against other algorithms such as CloSpan and SPADE. 2) Utilization of memory and momentum is rapid when compared to the CloSpan algorithm which is again analogous to BIDE. 3) There is the involvement of an enhanced occurrence and a probability reduction in the memory exploitation rate with the aid of the trait equivalent prognosis and also rim snipping of the PEPP with inference analysis for no coherent pattern pruning. This is on the basis of the surveillance done which concludes that PEPP's implementation is far more noteworthy and important in contrast with the likes of BIDE, to be precise.

JAVA 1.6_20th build was employed for accomplishment of the PEPP and BIDE algorithms. A workstation equipped with core2duo processor, 2GB RAM and Windows XP installation was made use of for investigation of the algorithms. The parallel replica was deployed to attain the thread concept in JAVA.

10 VII.

11 DATASET CHARACTERISTICS

Pi is supposedly found to be a very opaque dataset, which assists in excavating enormous quantity of recurring clogged series with a profitably high threshold somewhere close to 90%. It also has a distinct element of being enclosed with 190 protein series and 21 divergent objects. Reviewing of serviceable legacy's consistency has been made use of by this dataset. Fig. 5 portrays an image depicting dataset series extent status.

In assessment with all the other regularly quoted forms like SPADE, prefixspan and CloSpan, BIDE has made its mark as a most preferable, superior and sealed example of mining copy, taking in view the detailed study of the factors mainly, memory consumption and runtime, judging with PEPP. In contrast to PEPP and BIDE, a very intense dataset Pi is used which has petite recurrent closed series whose end to end distance is less than 10, even in the instance of high support amounting to around 90%. The diagrammatic representation displayed in Fig 3 explains that the above mentioned two algorithms execute in a similar fashion in case of support being 90% and above. But in situations when the support case is 88% and less, then the act of PEPP surpasses BIDE's routine. The disparity in memory exploitation of PEPP and BIDE can be clearly observed because of the consumption level of PEPP being lower than that of BIDE. The concept inference analysis we introduced here played a vital role in closed itemset detection. The significant improvement in closed itemset detection can be observable in our results, see the fig 6 VIII.

12 CONCLUSION

It has been scientifically and experimentally proved that clogged prototype mining propels dense product set and considerably enhanced competency as compared to recurrent prototype of mining even though both these types project similar animated power. The detailed study has verified that the case usually holds true when the count of recurrent moulds is considerably large and is the same with the recurrent bordered models as well. However, there is the downbeat in which the earlier formed clogged mining algorithms depend on chronological set of recurrent mining outlines. It is used to verify whether an innovative recurrent outline is blocked or else if it can nullify few previously mined blocked patterns. This leads to a situation where the memory utilization is considerably high but also leads to inadequacy of increasing seek out space for outline closure inspection. This paper anticipates an unusual algorithm for withdrawing recurring closed series with the help of Sequence Graph. It performs the following functions: It shuns the blight of contender's maintenance and test exemplar, supervises memory space expertly and ensures recurrent closure of clogging in a well-organized manner and at the same instant guzzling less amount of memory plot in comparison with the earlier developed mining algorithms. There is no necessity of preserving the already defined set of blocked recurrences, hence it very well balances the range of the count of frequent clogged models. A Sequence graph is embraced by PEPP and has the capability of harvesting the recurrent clogged pattern in an online approach. The efficacy of dataset drafts can be showcased by a wide-spread range of experimentation on a number of authentic datasets amassing varied allocation attributes. PEPP is rich in terms of velocity and memory spacing in comparison with the BIDE and CloSpan algorithms. ON the basis of the amount of progressions, linear scalability is provided. It is also proven that PEPP is efficient to find closed itemsets under inference analysis. It has been proven and verified by many scientific research studies that limitations are crucial for a number of chronological outlined mining algorithms. In addition we improved closed itemset detection performance by introducing inference analysis as an extension to PEPP. Future studies include proposing of post processing and pruning of the rules based on categorical relations between attributes.



Figure 1: 1 |

¹© 2011 Global Journals Inc. (US)

²© 2011 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XI Issue XIX Version I 2

³© 2011 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XI Issue XIX Version I

⁴November

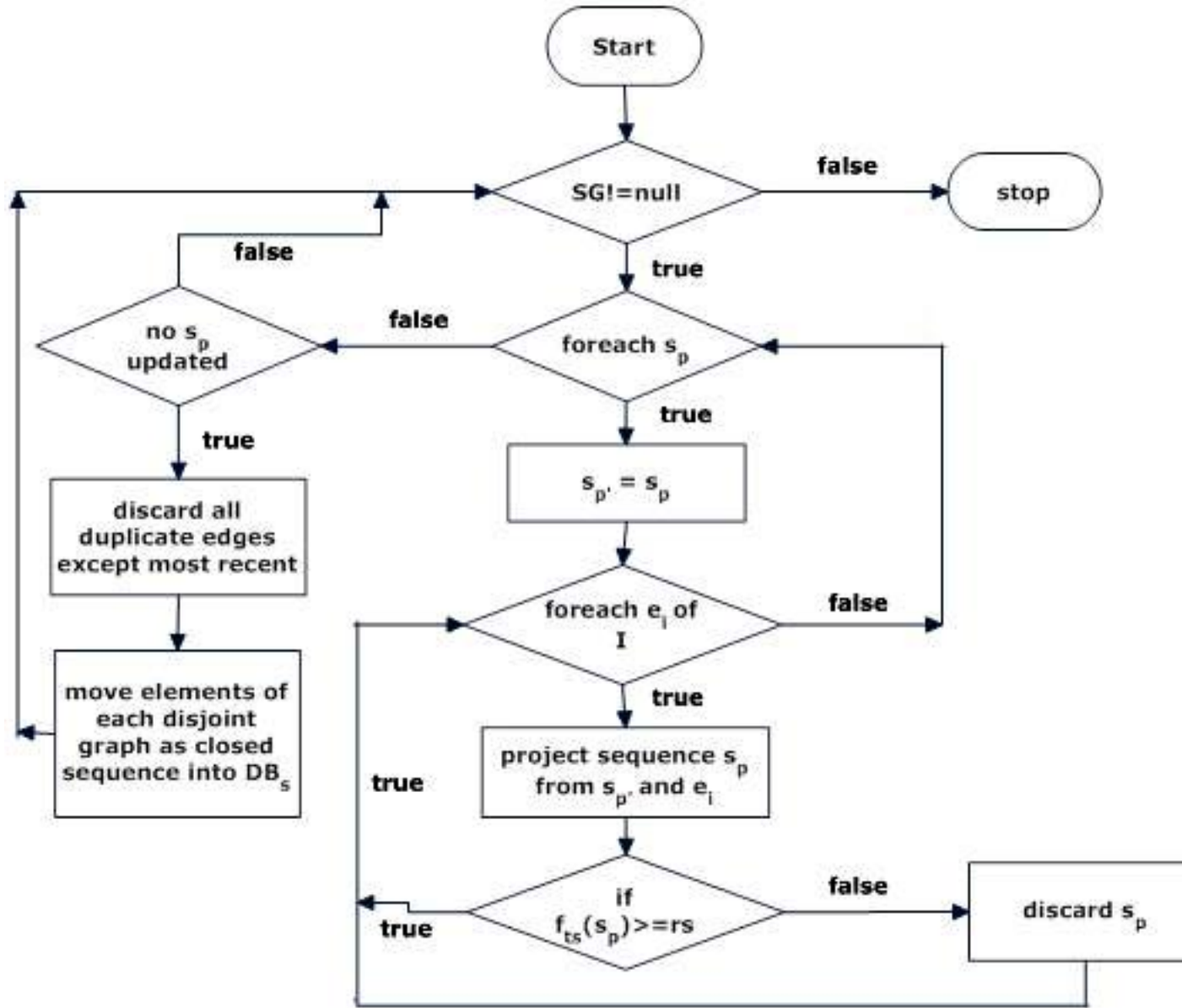
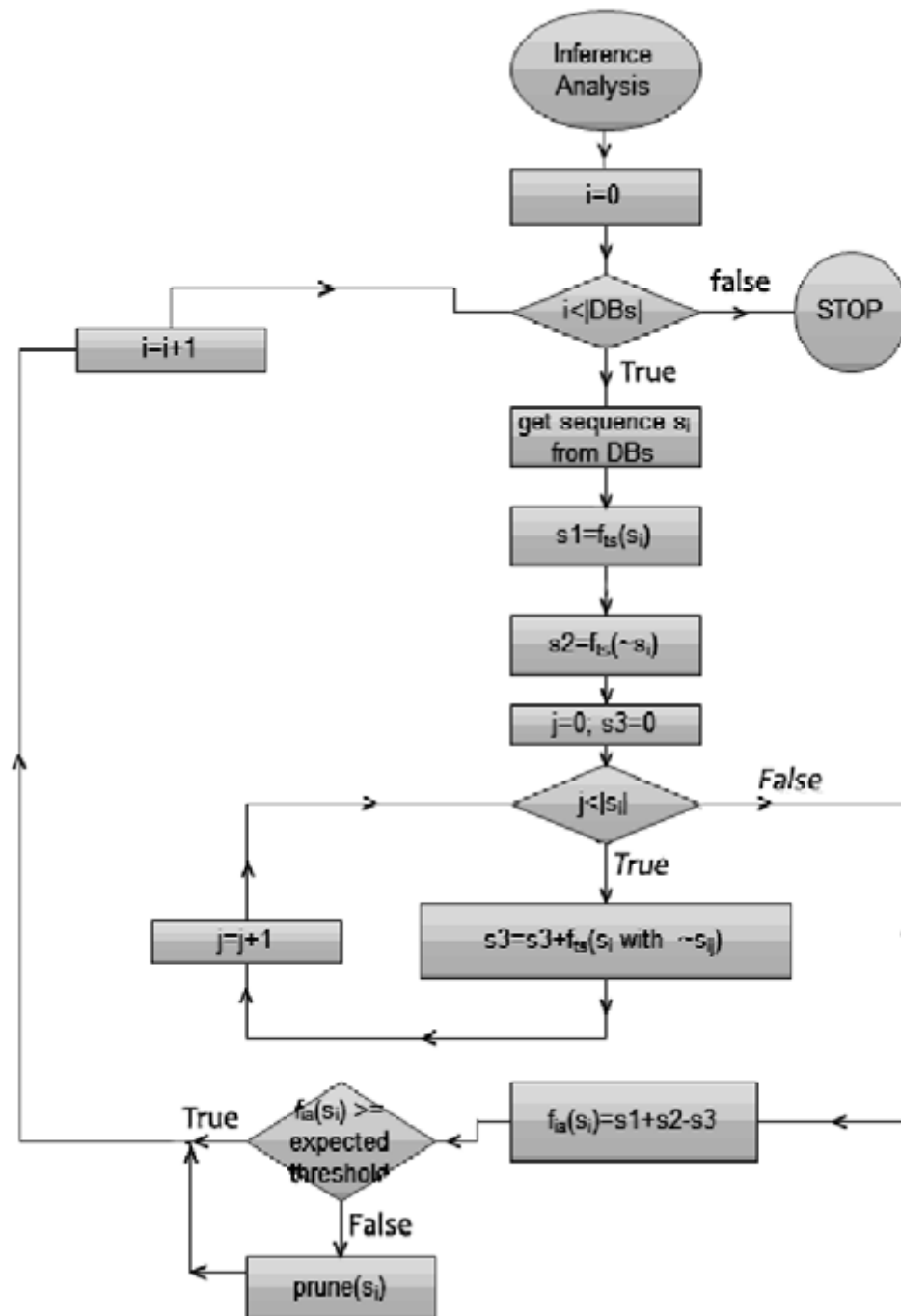


Figure 2:



2

Figure 3: Fig. 2 :

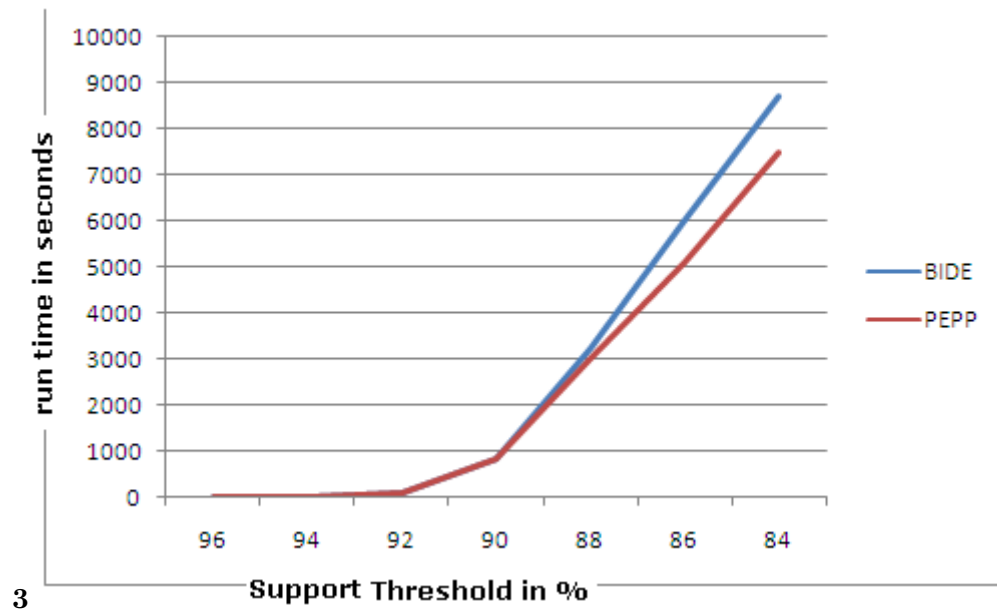


Figure 4: Fig. 3 :

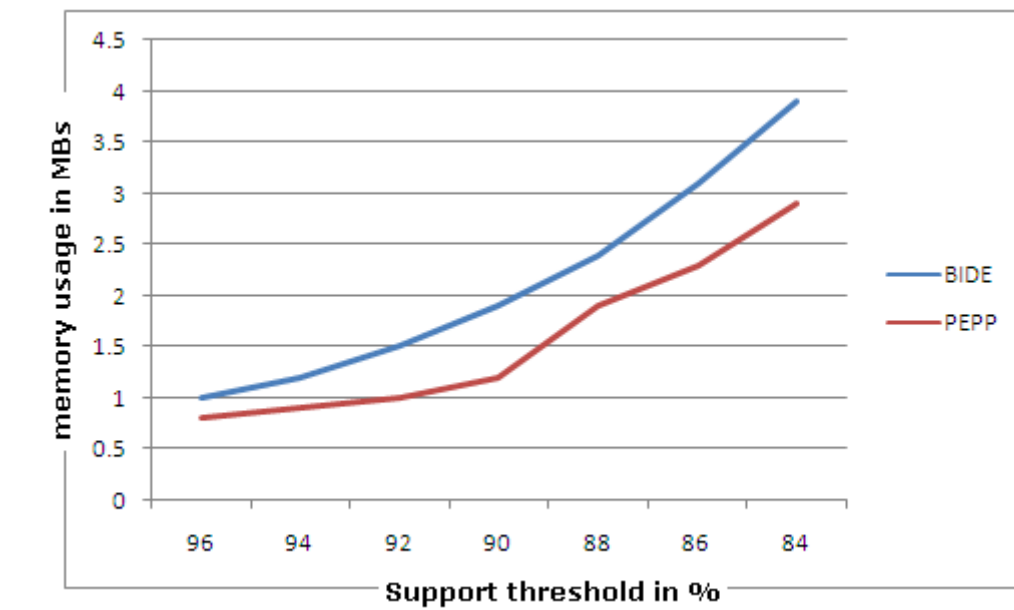


Figure 5: Fig. 4 :

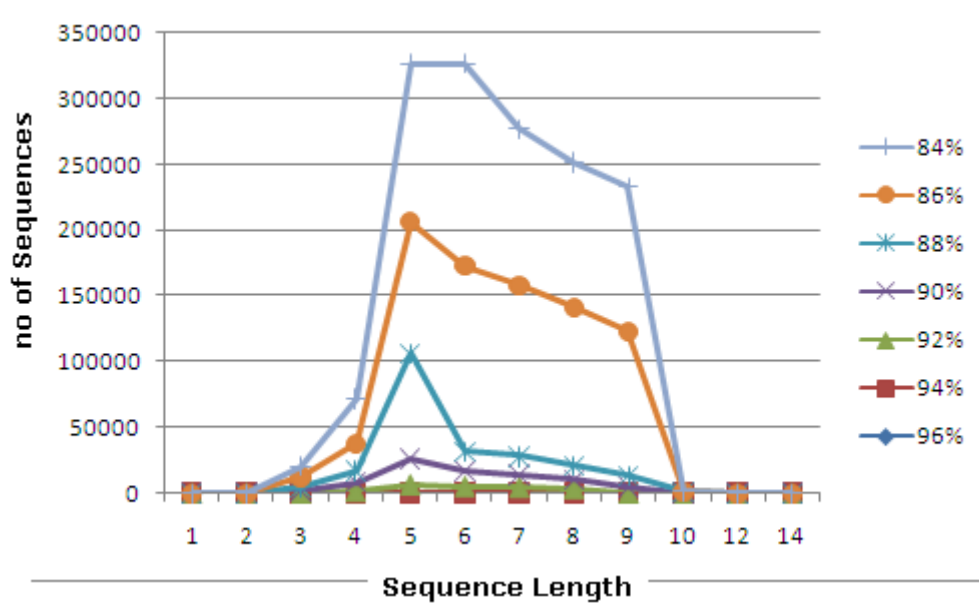


Figure 6: Fig. 5 :

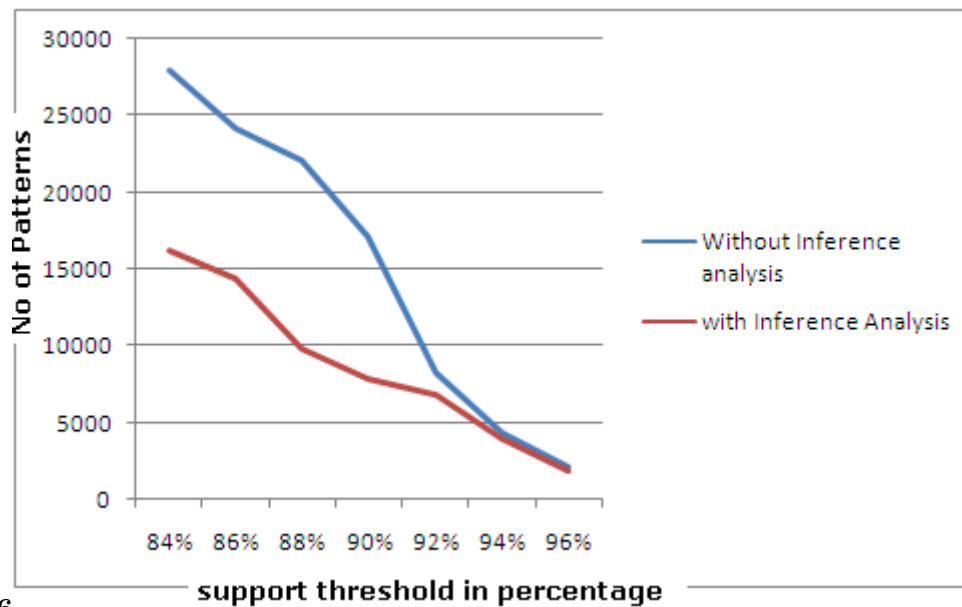


Figure 7: Fig. 6 :

$$\begin{aligned} & (s_i) \text{ ts } t \text{ f s } : (s_j) \text{ ts } t \text{ q s } t \text{ S f s f s DB} = \\ & \text{Sub-sequence} \end{aligned}$$

sequence is sub sequence for its next projected
sequence if both sequences having same total support. Super-sequence: A sequence is a super sequence for
N

super-sequence can be formulated as

If threshold given by user And $(s_i) \text{ ts } t \text{ f s } \geq r_s$ where 'rs' is required support : $t \text{ p s } s$ for any p value $<$
where

$$(s_i) \text{ t f s } \geq (s_j) \text{ p f s ts}$$

IV.

?

Total Support 'ts' : occur-
rence count of a
sequence as an ordered list
in all sequences in
sequence set 'S' can adopt as
total support 'ts' of that
sequence. Total support 'ts'
of a sequence can
determine by fallowing for-
mulation.
 $(s_i) \text{ t f s ts} = \sum (t \text{ p s } s)$
for each p

[Note: SDB is set of sequences.]

Figure 8:

12 CONCLUSION

L3: For each projected Itemset p s in memory
 Begin:

s $p' = s \setminus p$

L4: For each $i \in p'$

Begin:

Project $p \setminus s$ from C2: If $(p \setminus s) \cap p' \neq \emptyset$?
 (
 ,
)
 p
 i
 s
 e
 rs

Begin Spawn SG by adding edge between $p \setminus s$ and $e \setminus i$ End: C2 End: L4 C3: If p' is not spawned and no new

Remove all duplicate edges for each edge

weight from p'

Based on the above formulas we perform the logical analysis to derive the actual support of the patterns that improves the rule coherency. Inference analysis by example: Let A,B?I where I is itemset generated with the association of A,B are individual items or subsets.

Under logical analysis we determine ()

)

and ()

[Fayyad et al. ()] *Advances in Knowledge Discovery and Data Mining*, U M Fayyad , G Piatetsky-Shapiro , P Smyth , R Uthurusamy . 1996. AAAI/MIT Press.

[Kalli Srinivasa Nageswara Prasad and Prof. S Ramakrishna (2011)] ‘Article: Frequent Pattern Mining and Current State of the Art’. *International Journal of Computer Applications* Kalli Srinivasa Nageswara Prasad and Prof. S Ramakrishna (ed.) July 2011. 26 (7) p. . (Published by Foundation of Computer Science)

[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.

[Wang et al. (2003)] ‘CLOSET+: Searching for the Best Strategies for Mining Frequent Closed Itemsets’. J Wang , J Han , J Pei . *KDD’03*, (Washington, DC) Aug. 2003.

[Pei et al. (2001)] ‘CLOSET: An efficient algorithm for mining frequent closed itemsets’. J Pei , J Han , R Mao . *DMKD’01 workshop*, (Dallas, TX) May 2001.

[Yan et al. (2003)] ‘CloSpan: Mining Closed Sequential Patterns in Large Databases’. J Yan , R Han , Afshar . *SDM’03*, (San Francisco, CA) May 2003.

[Pei et al. (2002)] ‘Constraint-based sequential pattern mining in large databases’. J Pei , J Han , W Wang . *CIKM’02*, (McLean, VA) Nov. 2002.

[Ozden et al. (1998)] ‘Cyclic association rules’. B Ozden , S Ramaswamy , A Silberschatz . *ICDE’98*, (Orlando, FL) Feb. 1998.

[Pasquier et al. (1999)] ‘Discovering frequent closed itemsets for association rules’. N Pasquier , Y Bastide , R Taouil , L . *ICDT’99*, (Jerusalem, Israel) Jan. 1999.

[Mannila et al. (1995)] ‘Discovering frequent episodes in sequences’. H Mannila , H Toivonen , A I Verkamo . *SIGKDD’95*, (Montreal, Canada) Aug. 1995.

[Pasquier et al. ()] ‘Efficient Mining of Association Rules Using Closed Itemset Lattices’. N Pasquier , Y Bastide , R Taouil , L Lakhal . *Information Systems* 1999. 24 p. .

[Han et al. (1999)] ‘Efficient mining of partial periodic patterns in time series database’. J Han , G Dong , Y Yin . *ICDE’99*, (Sydney, Australia) Mar. 1999.

[Agrawal and Srikant (1994)] ‘Fast algorithms for mining association rules’. R Agrawal , R Srikant . *VLDB’94*, (Santiago, Chile) Sept. 1994.

[Jonassen et al. ()] ‘Finding flexible patterns in unaligned protein sequences’. J F Jonassen , D G Collins , Higgins . *Protein Science* 1995. 4 (8) .

[Han et al. (2000)] ‘FreeSpan: Frequent patternprojected sequential pattern mining’. J Han , J Pei , B Mortazavi-Asl , Q Chen , U Dayal , M C Hsu . *SIGKDD’00*, (Boston, MA) Aug. 2000.

[Zaki ()] ‘Generating Non-Redundant Association Rules’. M J Zaki . *Proc. Int’l Conf. Knowledge Discovery and Data Mining*, (Int’l Conf. Knowledge Discovery and Data Mining) 2000. p. .

[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.

[Machine Learning, and Knowledge Discovery in Databases ()] *Machine Learning, and Knowledge Discovery in Databases*, 1995. p. . (Workshop Statistics)

[Burdick et al. (2005)] ‘Mafia: A Maximal Frequent Itemset Algorithm’. D Burdick , M Calimlim , J Flannick , J Gehrke , T Yiu . *IEEE Trans. Knowledge and Data Eng* Nov. 2005. 17 (11) p. .

[Agrawal et al. ()] ‘Mining Association Rules between Sets of Items in Large Databases’. R Agrawal , T Imielinski , A Swami . *Proc. ACM SIGMOD*, (ACM SIGMOD) 1993. p. .

[Yang et al. (2002)] ‘Mining long sequential patterns in a noisy environment’. J Yang , P S Yu , W Wang , J Han . *SIGMOD’02*, (Madison, WI) June 2002.

[Agrawal and Srikant (1995)] ‘Mining sequential patterns’. R Agrawal , R Srikant . *ICDE’95*, (Taipei, Taiwan) Mar. 1995.

- [Srikant (1996)] ‘Mining sequential patterns: Generalizations and performance improvements’. R Srikant , R .
EDBT’96, (Avignon, France) Mar. 1996.
- [Bettini et al. ()] ‘Mining temporal relationals with multiple granularities in time sequences’. C Bettini , X Wang
, S Jajodia . *Data Engineering Bulletin* 1998. 21 (1) 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.
- [Li (2006)] ‘On Optimal Rule Discovery’. J Li . *IEEE Trans. Knowledge and Data Eng* Apr. 2006. 18 (4) p. .
- [Baesens et al. ()] ‘Post-Processing of Association Rules’. B Baesens , S Viaene , J Vanthienen . *Proc. Workshop
Post-Processing in Machine Learning and Data Mining: Interpretation, Visualization, Integration, and
Related Topics with Sixth ACM SIGKDD*, (Workshop Post-essing in Machine Learning and Data Mining:
Interpretation, Visualization, Integration, and Related Topics with Sixth ACM SIGKDD) 2000. p. .
- [Pei et al. (2001)] ‘Prefixspan: Mining sequential patterns efficiently by prefix-projected pattern growth’. J Pei ,
J Han , B Mortazavi-Asl , Q Chen , U Dayal , M C Hsu . *ICDE’01*, (Heidelberg, Germany) April 2001.
- [Toivonen et al.] ‘Pruning and Grouping of Discovered Association Rules’. H Toivonen , M Klemettinen , P
Ronkainen , K Hatonen , H Mannila . *Proc. ECML-95*, (ECML-95)
- [Ayres et al. (2002)] ‘Sequential Pattern Mining using a Bitmap Representation’. J Ayres , J Gehrke , T Yiu , J
Flannick . *SIGKDD’02*, (Edmonton, Canada) July 2002.
- [Seno and Karypis (2002)] ‘SLPMiner: An algorithm for finding frequent sequential patterns using length
decreasing support constraint’. M Seno , G Karypis . *ICDM’02*, (Maebashi, Japan) Dec. 2002.
- [Zaki ()] ‘SPADE: An Efficient Algorithm for Mining Frequent Sequences’. M Zaki . *Machine Learning*, 2001.
Kluwer Academic Publishers. 42 p. .
- [Garofalakis et al. (1999)] ‘SPIRIT: Sequential Pattern Mining with regular expression constraints’. M Garo-
falakis , R Rastogi , K Shim . *VLDB’99*, (San Francisco, CA) Sept. 1999.
- [Massegia et al. (1995)] ‘The psp approach for mining sequential patterns’. F Massegia , F Cathala , P Poncelet
. *PKDD’98*, (Nantes, France) Sept. 1995.
- [Zaki and Ogihara (1998)] ‘Theoretical Foundations of Association Rules’. M J Zaki , M Ogihara . *Proc. Workshop
Research Issues in Data Mining and Knowledge Discovery (DMKD ’98)*, (Workshop Research Issues in Data
Mining and Knowledge Discovery (DMKD ’98)) June 1998. p. .
- [Silberschatz and Tuzhilin (1996)] ‘What Makes Patterns Interesting in Knowledge Discovery Systems’. A
Silberschatz , A Tuzhilin . *IEEE Trans. Knowledge and Data Eng* Dec. 1996. 8 (6) p. .