

# CMRULES: An Efficient Algorithm for Mining Sequential Rules Common to Several Sequences

Philippe Fournier-Viger<sup>1</sup>

Usef Faghihi<sup>1</sup>

Roger Nkambou<sup>1</sup>

Engelbert Mephu Nguifo<sup>2</sup>

<sup>1</sup>University of Québec in Montréal (Canada)

<sup>2</sup>University Blaise-Pascal(France)



May, 20<sup>th</sup> 2010

## Introduction

- **Discovering temporal relationships** between events stored in large databases is important in many domains. It helps to understand the relationships between events and sets a basis for the prediction of events
- For example:
  - predicting stock market
  - predicting consumer behavior
- In this paper, we focus on temporal relationships in sequences of discrete events

## Sequential Pattern Mining

- One of the most popular set of techniques for discovering temporal relations between events in sequences of event is **sequential pattern mining**.
- It finds subsequences that are common to several sequences in a sequence database.
- However, knowing that a sequence of events appear frequently in a database **is not sufficient for the prediction of events**. For example, it is possible that some event  $y$  appears frequently after an event  $x$  but that there are also many cases where  $x$  is not followed by  $y$ .
- For **prediction**, we need a measurement of the confidence that if  $x$  occurs  $y$  will occur afterward.

3

## Sequential Rule Mining

- The alternative that addresses prediction is **Sequential Rule Mining**.
- **Sequential rules** are also called episode rules, prediction rule or temporal rules.
- A sequential rule typically has a *confidence* and a *support*.
- Manila et al. (1997) discovers rules of the form  $X \rightarrow Y$  appearing frequently **in a single sequence of events**.  $X$  and  $Y$  are sets of events.
- Other works that find rules **in a single sequence of events**: Hamilton & Karimi (2005), Hsieh (2006), Deogun (2005), etc.
- **Several applications**: stock market analysis (Das et al., 1998; Hsieh et al., 2006), weather observation (Hamilton & Karimi, 2005), drought management (Harms et al. 2002), alarm analysis, etc.

4

## Sequential Rule Mining (continued)

- Algorithms for finding rules occurring **frequently in several sequences**:  
Das et al. (1998),  $x \rightarrow y$  where  $x, y$  are single events  
Harm et al. (2002).  $X \rightarrow Y$  where  $X, Y$  are sets of events
- However, these algorithms are **not** appropriate for finding **rules that are common to several sequences**.
- The reason is that they are designed for mining rules appearing frequently in sequences. With these algorithms, if a rule appears many times in a sequence, it can be considered frequent even if it does not appear in any other sequences.
- In this paper, **we address the problem of mining sequential rule common to several sequences**.

5

## Our proposal

- We propose an algorithm for mining rules **common to several sequences**. This can be useful for example for mining rules that are common to several customers.
- The algorithm is based on association rule mining.
- The idea is to find associations rules between items to prune the search space to items that occur jointly in many sequences. Then it eliminates association rules that are not sequential rules. We prove that this correctly generates all sequential rules.

6

## Association Rule Mining

- A **transaction database**  $D$  is a set of transactions  $T=\{t_1, t_2, \dots, t_n\}$  and a set of items  $I=\{i_1, i_2, \dots, i_n\}$ , where  $t_1, t_2, \dots, t_n \subseteq I$ .
- **Association rule:**  
 $X \rightarrow Y$ , such that  $X, Y \subseteq I, X \cap Y = \emptyset$ ,
- The **support** of a rule  $X \rightarrow Y$  is defined as  $\text{sup}(X \cup Y) / |T|$ . The **confidence** of a rule is defined as  $\text{conf}(X \rightarrow Y) = \text{sup}(X \cup Y) / \text{sup}(X)$ .
- **Association rule mining:** finding all rules such that their support and confidence are higher to thresholds *minsup* and *minconf*.

ID	transactions
1	abcefg
2	abcdef
3	abef
4	bfg

(2) Association rule mining

↓  
Association rules

ID	rule	Sup.	Conf.
1	$c \rightarrow abef$	0.5	1.0
2	$abc \rightarrow e$	0.5	1.0
3	$a \rightarrow cef$	0.5	0.6
4	$ab \rightarrow efg$	0.75	1.0
5	$b \rightarrow efg$	0.75	0.75
6	$a \rightarrow efg$	0.75	1.0
7	$c \rightarrow f$	0.5	1.0
8	$a \rightarrow b$	0.75	1.0
...	...	...	...

7

## Definition of a Sequence Database

- A **sequence database**  $SD$  is defined as a set of sequences  $SD=\{s_1, s_2, \dots, s_n\}$  and a set of items  $I=\{i_1, i_2, \dots, i_n\}$ , where each sequence  $s_x$  is an ordered list of transactions  $s_x=\{X_1, X_2, \dots, X_n\}$  such that  $X_1, X_2, \dots, X_n \subseteq I$ .
- An example:

ID	sequences
1	(ab), (c), (f), (g), (e)
2	(ad), (c), (b), (ef)
3	(a), (b), (f) (e)
4	(b), (fg)

8

## Mining Sequential Rules Common to Several Sequences

- A **sequential rule**  $X \Rightarrow Y$  is a relationship between two itemsets  $X, Y$  such that  $X, Y \subseteq I$  and  $X \cap Y = \emptyset$ .
- **Sequential support:**  $\text{seqSup}(X \Rightarrow Y) = \text{sup}(X \blacksquare Y) / |SD|$ .
- **Sequential confidence:**  $\text{seqConf}(X \Rightarrow Y) = \text{sup}(X \blacksquare Y) / \text{sup}(X)$ .
- $\text{sup}(X \blacksquare Y)$  denotes the number of sequences from a sequence database where all the items of  $X$  appear before all the items of  $Y$ .
- **Sequential rule mining:** finding all rules with sequential support and a sequential confidence higher than some thresholds  $\text{minSeqSup}$  and  $\text{minSeqConf}$ .

9

## An example of Sequential Rule Mining

Some rules found for  $\text{minSeqSup} = 0.5$  and  $\text{minSeqConf} = 0.5$ :

ID	sequences	ID	rule	S-Sup.	S-Conf.
1	(a b), (c), (f), (g), (e)	1	a b c $\Rightarrow$ e	0.5	1.0
2	(a d), (c), (b), (e f)	2	a $\Rightarrow$ c e f	0.5	0.66
3	(a), (b), (f) (e)	3	a b $\Rightarrow$ e f	0.5	1.0
4	(b), (f g)	4	b $\Rightarrow$ e f	0.75	0.75
		5	a $\Rightarrow$ e f	0.75	1.0
		6	c $\Rightarrow$ f	0.5	1.0
		7	a $\Rightarrow$ b	0.5	0.66
		...	...	...	...

10

## The Relationship between Association Rules and Sequential Rules

- A sequence database **S** can be transformed into a transaction database **S'** by removing time information.
- Each sequential rule  $r: X \Rightarrow Y$  of **S** has a corresponding association rule  $r': X \rightarrow Y$  in **S'**.
- Since  $\text{sup}(X \blacksquare Y)$  is always lower or equal to  $\text{sup}(X \cup Y)$  these relationships hold:
  - $\text{sup}(r') \geq \text{seqSup}(r)$
  - $\text{conf}(r') \geq \text{seqConf}(r)$ ,for any sequential rule  $r$  and its corresponding association rule  $r'$ .

11

## The CMRules algorithm

**INPUT** : a sequence database,  $\text{minSeqSup}$ ,  $\text{minSeqConf}$

**OUTPUT** : the set of all sequential rules

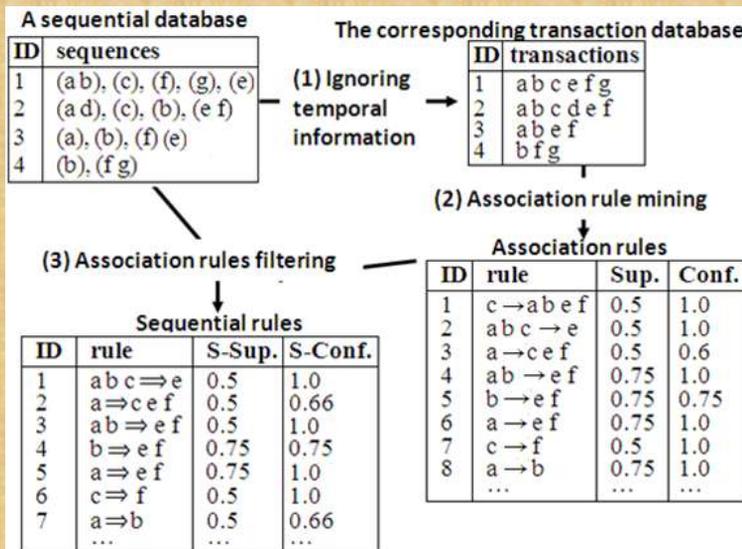
**PROCEDURE**:

1. Consider the sequence database as a transaction database
2. Find all association rules from the transaction database by applying an association rule mining algorithm such as Apriori (Agrawal et al., 1993). Select  $\text{minsup} = \text{minSeqSup}$  and  $\text{minconf} = \text{minSeqConf}$ .
3. Scan the original sequence database to calculate the sequential support and sequential confidence of each association rule found in the previous step. Eliminate each rule  $r$  such that:
  - a.  $\text{seqSup}(r) < \text{minSeqSup}$
  - b.  $\text{seqConf}(r) < \text{minSeqConf}$
4. Return the set of rules

see the article for the proof of correctness.

12

## A Sample Execution of the CMRules algorithm



13

## Optimizations

- How to calculate  $\text{sup}(X \blacksquare Y)$  for a rule  $X \rightarrow Y$  ?
  - The naïve approach is to check all sequences.
  - Association rule mining algorithms first find frequent itemsets  $X$  and  $Y$  and then generate rules of the form:  $X \rightarrow Y - X$
  - Algorithms such as AprioriTID, Eclat, H-Mine, etc. can annotate itemset  $X$  and  $Y$  with the sequences that contains them.
  - If we use such algorithm, we can only check sequence containing  $X$  for calculating  $\text{sup}(X \blacksquare Y)$ .
  - This can improve performance by up to 50 %.
- Also, we don't need to keep the association rules. We can discard each rule immediately after it is found if it is not a sequential rule. This reduces memory consumption.

14

## Analysis of the Time Complexity

- Step 1: converting a sequence database in a transaction database is done in linear time.
- Step 2: association rule mining. Two substeps:
  1. Discovering frequent itemsets:  
Apriori :  $O(d^2 n)$   $d$ = number of diff. items,  $n$ = number of transactions
  2. Generating association rules from frequent itemsets:  
less costly, thus can be ignored.
- Step 3: calculating sequential conf. and support.  
For each rule, **best case**:  $|S| \times \text{minsup}$  sequences to check, **worst case**:  $|S|$  sequences to check. Checking if a rule is included in a sequence is done in linear time.

15

## A Second Algorithm: CMDeogun

- It is an adaptation of the algorithm of Deogun et al. for the case of several sequences.
- The approach is similar to the Apriori algorithm.
- The process:
  - **First** find all rules with an antecedent and consequent of size 1.
  - **Then**, recursively combine rules of smaller size to discover larger rules by performing left expansion and right expansion.

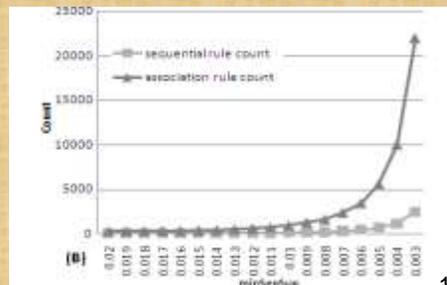
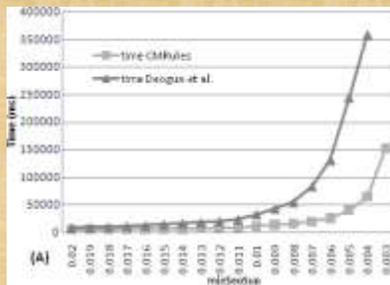
16

## Comparison of CMRules and CMDeogun on the Kosarak dataset

**Kosarak Dataset** : 990 000 click-stream data from the logs of online news portal.

**Experiment**, we selected the 70 000 first sequences of the dataset. These sequences contain an average of 7.97 items each, from a set of 21 144 different items.

We applied the algorithms on this dataset with  $minSeqConf = 0.3$  and with  $minSeqSup = 0.02, 0.019, \dots, 0.003$ .



17

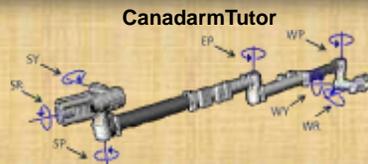
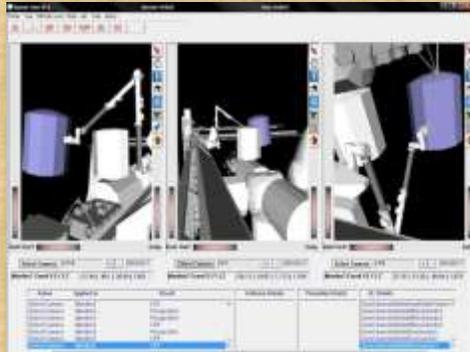
## Discussion of results

- The percentage of sequential rules can be quite high in real datasets.
  - For  $minSeqSup$  of 0.02, there is 312 association rules and 64 sequential rules (20.51 %).
  - For a  $minSeqSup$  of 0.003, there is 21 947 association rules and 2 512 sequential rules (11.44%).
- The execution time of CMRules seems to grow proportionally to the number of association rules.
- CMRules has a better scalability for this dataset.
- The time for converting the sequence database in a transaction database is negligible (~300 ms).

18

## An Application of CMRules in a Cognitive Agent

- The **CTS Cognitive Agent** has some predefined behaviors.
- Each execution of CTS is recorded as a sequence in a sequence database
- Items = actions and perceptions of CTS
- Itemset = one cognitive cycle
- Rules are used for making predictions about which behavior will be successful with learners, and adapt the behavior of CTS accordingly.



19

## An Application of CMRules in a Cognitive Agent

- We did 100 executions of CTS (100 sequences)
- Sequence length: average of 25 itemsets
- CMRules was executed after each execution of CTS
- On average 31.05% of the association rules were sequential rules.
- Execution time was less than 50 ms.
- We have found several relevant rules. For example:
  - “the learner move the wrong joint  $\Rightarrow$  the learner makes the robotic arm pass too close to the space station”
  - “the learner lacks motivation  $\Rightarrow$  the learner is inactive”

20

## Conclusion

### Summary

- CMRules: an algorithm for mining sequential rules common to several sequences.
- CMRules can discover sequential rules and association rules at the same time
- Extensions by using another association rule mining algorithm (incremental mining, parallel mining, etc.).
- Extensions by modifying the definition of  $\text{sup}(X \blacksquare Y)$
- CMDeogun: a second algorithm.

### Future works

- Further experiments will be done to measure empirically the performance of CMRules and CMDeogun. We will also try to define other algorithms.

21

## Thank you. Questions?



Thanks to the organizers of the FLAIRS 2010 conference!

### Acknowledgement to:

.Current and past members of the GDAC and PLANIART research teams for their collaboration.  
.NSERC and FQRNT funding programs

Fonds de recherche  
sur la nature  
et les technologies  
Québec

