

WEB USAGE MINING: EXTRACTION OF SEQUENTIAL PATTERNS

USING PREFIXSPAN METHOD

SAHANA V¹ & G T RAJU²

¹Dean of Enggining, Department of Computer Science, RNS Institute of Technology, Bangalore, Karnataka, India

²Professor & H.O.D, Department of Computer Science, RNS Institute of Technology, Bangalore, Karnataka, India

ABSTRACT

Sequential pattern mining is an important data mining problem with broad applications. However, it is also a difficult problem since the mining may have to degenerate or examine a combinatorially explosive number of intermediate subsequences. Most of the previously developed sequential pattern mining methods, such as GSP, explore a candidate generation-and-test approach [1] to reduce the number of candidates to be examined. However, this approach may not be efficient in mining large sequence databases having numerous patterns and/or long patterns. In this paper, we propose a projection-based, sequential pattern-growth approach for efficient mining of sequential patterns. In this approach, a sequence database is recursively projected into a set of smaller projected databases, and sequential patterns are grown in each projected database by exploring only locally frequent fragments.

KEYWORDS: Sequential Pattern, Frequent Pattern, Transaction Database, Sequence Database

I. INTRODUCTION

Sequential pattern mining, which discovers frequent subsequences as patterns in a sequence database, is an important data mining problem with broad applications, including the analysis of customer purchase patterns or Web access patterns, the analysis of sequencing or timerelated processes such as scientific experiments, natural disasters, and disease treatments, the analysis of DNA sequences, etc.

The sequential pattern mining problem was first introduced by Agrawal and Srikant in [2]: Given a set of sequences, where each sequence consists of a list of elements and each element consists of a set of items, and given a user-specified min_support threshold, sequential pattern mining is to find all frequent subsequences, i.e., the subsequences whose occurrence frequency in the set of sequences is no less than min_support.

The a priori-like sequential pattern mining method, though reducing search space, bears three nontrivial, inherent costs that are independent of detailed implementation techniques.

- A huge set of candidate sequences could be generated in a large sequence database. Since the set of candidate sequences includes all the possible permutations of the elements and repetition of items in a sequence, the a priori-based method may generate a really large set of candidate sequences even for a moderate seed set. For example, two frequent sequences of length-1, $\langle a \rangle$ and $\langle b \rangle$ will generate five candidate sequences of length-2: $\langle aa \rangle$, $\langle ab \rangle$, $\langle ba \rangle$, $\langle bb \rangle$, and $\langle (ab) \rangle$, where $\langle (ab) \rangle$ represents that two events, a and b, happen in the same time slot. If there are 1,000 frequent sequences of length-1, such as $\langle a_1 \rangle$, $\langle a_2 \rangle$, $\langle a_{1000} \rangle$, an a priori-like algorithm will generate $1,000 * 1,000 + (1,000 + 999) / 2 = 1,499,500$ candidate sequences. Notice that the cost of candidate sequence

generation, test, and support counting is inherent to the a priori-based method, no matter what technique is applied to optimize its detailed implementation.

- Multiple scans of databases in mining. The length of each candidate sequence grows by one at each database scan. In general, to find a sequential pattern of length l , the a priori-based method must scan the database at least l times. This bears a nontrivial cost when long patterns exist.
- The a priori-based method generates a combinatorially explosive number of candidates when mining long sequential patterns. A long sequential pattern contains a combinatorial explosive number of subsequences, and such subsequences must be generated and tested in the a priori-based mining. Thus, the number of candidate sequences is exponential to the length of the sequential patterns to be mined. For example, let the database contain only one single sequence of length 100, $\langle a_1 a_2 \dots a_{100} \rangle$, and the min_support threshold be 1 (i.e., every occurring pattern is frequent). To (re)derive this length-100 sequential pattern, the a priori-based method has to generate 100 length-1 candidate sequences (i.e., $\langle a_1 \rangle, \langle a_2 \rangle, \dots, \langle a_{100} \rangle$), $100 * 100 + (100 + 99) / 2 = 14,950$ length-2 candidate sequences, $\binom{100}{3} = 161,700$ length-3 candidate sequences, and so on. Obviously, the total number of candidate sequences to be generated is

$$\sum_{i=1}^{100} \binom{100}{i} = 2^{100} - 1 \approx 10^3.$$

In this paper, we systematically explore a pattern-growth approach for efficient mining of sequential patterns in large sequence database. The approach adopts a divide-and-conquer, pattern-growth principle as follows: Sequence databases are recursively projected into a set of smaller projected databases based on the current sequential pattern(s), and sequential patterns are grown in each projected databases by exploring only locally frequent fragments.

Table 1: A Sequence Database

Sequence_id	Sequence
10	$\langle a(abc)(ac)d(cf) \rangle$
20	$\langle (ad)c(bc)(ae) \rangle$
30	$\langle (ef)(ab)(df)cb \rangle$
40	$\langle eg(af)cbc \rangle$

The remainder of the paper is organized as follows: In Section 2, the sequential pattern mining problem is defined. In Section 3, our approach, projection based sequential pattern growth, is introduced. In section 4, our experimental results are reported.

II. PROBLEM DEFINITION

Let $I = \{i_1, i_2, \dots, i_n\}$ be a set of all items. An itemset is a subset of items. A sequence is an ordered list of item sets. A sequence s is denoted $\langle s_1 s_2 \dots s_l \rangle$, where s_j is an itemset. s_j is also called an element of the sequence, and denoted as $(x_1 x_2 \dots x_m)$, where x_k is an item. For brevity, the brackets are omitted if an element has only one item, i.e., element (x) is written as x . An item can occur at most once in an element of a sequence, but can occur multiple times in different elements of a sequence. The number of instances of items in a sequence is called the length of the sequence. A sequence with length l is called an l -sequence. A sequence $\alpha = \langle a_1 a_2 \dots a_n \rangle$ is called a subsequence of another sequence $\beta = \langle b_1 b_2 \dots b_m \rangle$ and

β a supersequence of α , denoted as $\alpha \sqsubseteq \beta$, if there exist integers $1 \leq j_1 < j_2 < \dots < j_n \leq m$ such that $a_1 \sqsubseteq b_{j_1}, a_2 \sqsubseteq b_{j_2}, \dots, a_n \sqsubseteq b_{j_n}$.

A sequence database S is a set of tuples $\langle sid, s \rangle$, where sid is a sequence_id and s a sequence. A tuple $\langle sid, s \rangle$ is said to contain sequence α , if α is a subsequence of s . The support of a sequence α in a sequence database S is the number of tuples in the database containing α , i.e., $support_S(\alpha) = |\{ \langle sid, s \rangle \mid (\langle sid, s \rangle \in S) \wedge (\alpha \sqsubseteq s) \}|$. It can be denoted as $support(\alpha)$ if the sequence database is clear from the context. Given a positive integer $min_support$ as the support threshold, a sequence α is called a sequential pattern in sequence database S if $support_S(\alpha) \geq min_support$. A sequential pattern with length l is called l -pattern.

Given a sequence database and the $min_support$ threshold, sequential pattern mining is to find the complete set of sequential patterns in the database.

To avoid checking every possible combination of a potential candidate sequence, one can first fix the order of items within each element. Since items within an element of a sequence can be listed in any order, without loss of generality, one can assume that they are always listed alphabetically. For example, the sequence in S with Sequence_id 10 in our running example is listed as $\langle a(abc)(ac)d(cf) \rangle$ instead of $\langle a(bac)(ca)d(fc) \rangle$. With such a convention, the expression of a sequence is unique. Then, the task is to examine whether one can fix the order of item projection in the generation of a projected database. Intuitively, if one follows the order of the prefix of a sequence and projects only the suffix of a sequence, one can examine in an orderly manner all the possible subsequences and their associated projected database. Thus, it is important to introduce the concept of prefix, suffix and projected database.

III. MINING SEQUENTIAL PATTERNS BY PATTERN GROWTH

Instead of repeatedly scanning the entire database and generating and testing large sets of candidate sequences, one can recursively project a sequence database into a set of smaller databases associated with the set of patterns mined so far and, then, mine locally frequent patterns in each projected database.

PrefixSpan: Prefix-Projected Sequential

Patterns Mining

We examine whether one can fix the order of item projection in the generation of a projected database. Intuitively, if one follows the order of the prefix of a sequence and projects only the suffix of a sequence, one can examine in an orderly manner all the possible subsequences and their associated projected database. Thus, we first introduce the concept of prefix and suffix.

Definition 1 (Prefix)

Suppose all the items within an element are listed alphabetically. Given a sequence

$\alpha = \langle e_1 e_2 \dots e_n \rangle$ (where each e_i corresponds to a frequent element in S), a sequence $\beta = \langle e'_1 e'_2 \dots e'_m \rangle$ ($m \leq n$) is called a prefix of α if and only if

- $e'_i = e_i$ for $(i \leq m - 1)$
- $e'_m \sqsubseteq e_m$

- all the frequent items in $(e_m - e'_m)$ are alphabetically after those in e'_m

For example, $\langle a \rangle, \langle aa \rangle, \langle a(ab) \rangle$ and $\langle a(abc) \rangle$ are prefixes of sequence $s(a(abc)(ac)d(cf))$, but neither $\langle ab \rangle$ nor $\langle a(bc) \rangle$ is considered as a prefix if every item in the prefix $\langle a(abc) \rangle$ of sequence s is frequent in S .

Definition 2 (Suffix)

Given a sequence $\alpha = \langle e_1 e_2 \dots e_n \rangle$ (where each e_i corresponds to a frequent element in S). Let $\beta = \langle e_1 e_2 \dots e_{m-1} e'_m \rangle (m \leq n)$ be the prefix of α . Sequence $\gamma = \langle e''_m e_{m+1} \dots e_n \rangle$ is called the suffix of α with regards to prefix β , denoted as $\gamma = \alpha / \beta$, where $e''_m = (e_m - e'_m)$. We also denote $\alpha = \beta . \gamma$.

For example, for the sequence $s = \langle a(abc)(ac)d(cf) \rangle, \langle (abc)(ac)d(cf) \rangle$ is the suffix with regards to the prefix $\langle a \rangle, \langle (_bc)(ac)d(cf) \rangle$ is the suffix with regards to the prefix $\langle aa \rangle$, and $\langle (_c(ac)d(cf)) \rangle$ is the suffix with regards to the prefix $\langle a(ab) \rangle$.

Definition 3(Projected Database)

Let α be a sequential pattern in a sequence database S . The α -projected database, denoted as $S|_\alpha$, is the collection of suffixes of sequences in S with regards to prefix α .

For the same sequence database S in Table 1 with $\text{min_sup} = 2$, sequential patterns in S can be mined by a prefix-projection method in the following steps:

Find Length-1 Sequential Patterns. Scan S once to

Find all the frequent items in sequences. Each of these frequent items is a length-1 sequential pattern.

They are $\langle a \rangle : 4, \langle b \rangle : 4, \langle c \rangle : 4, \langle d \rangle : 3, \langle e \rangle : 3$, and

$\langle f \rangle : 3$, where the notation “ $\langle \text{pattern} \rangle : \text{count}$ ” represents

The pattern and its associated support count.

Table 2: Projected Databases and Sequential Patterns

prefix	projected (suffix) database	sequential patterns
$\langle a \rangle$	$\langle (abc)(ac)d(cf) \rangle, \langle (_d)c(bc)(ae) \rangle, \langle (_b)(df)cb \rangle, \langle (_f)cbc \rangle$	$\langle a \rangle, \langle aa \rangle, \langle ab \rangle, \langle a(bc) \rangle, \langle a(bc)a \rangle, \langle aba \rangle, \langle abc \rangle, \langle (ab) \rangle, \langle (ab)c \rangle, \langle (ab)d \rangle, \langle (ab)f \rangle, \langle (ab)dc \rangle, \langle ac \rangle, \langle aca \rangle, \langle acb \rangle, \langle acc \rangle, \langle ad \rangle, \langle adc \rangle, \langle af \rangle$
$\langle b \rangle$	$\langle (_c)(ac)d(cf) \rangle, \langle (_c)(ae) \rangle, \langle (df)cb \rangle, \langle c \rangle$	$\langle b \rangle, \langle ba \rangle, \langle bc \rangle, \langle (bc) \rangle, \langle (bc)a \rangle, \langle bd \rangle, \langle bdc \rangle, \langle bf \rangle$
$\langle c \rangle$	$\langle (ac)d(cf) \rangle, \langle (bc)(ae) \rangle, \langle b \rangle, \langle bc \rangle$	$\langle c \rangle, \langle ca \rangle, \langle cb \rangle, \langle cc \rangle$
$\langle d \rangle$	$\langle (cf) \rangle, \langle c(bc)(ae) \rangle, \langle (_f)cb \rangle$	$\langle d \rangle, \langle db \rangle, \langle dc \rangle, \langle dcb \rangle$
$\langle e \rangle$	$\langle (_f)(ab)(df)cb \rangle, \langle (af)cbc \rangle$	$\langle e \rangle, \langle ea \rangle, \langle eab \rangle, \langle eac \rangle, \langle each \rangle, \langle eb \rangle, \langle ebc \rangle, \langle ec \rangle, \langle ecb \rangle, \langle ef \rangle, \langle efb \rangle, \langle efc \rangle, \langle efc b \rangle$
$\langle f \rangle$	$\langle (ab)(df)cb \rangle, \langle cbc \rangle$	$\langle f \rangle, \langle fb \rangle, \langle fbc \rangle, \langle fc \rangle, \langle fcb \rangle$

Divide Search Space

The complete set of sequential patterns can be partitioned into the following six subsets according to the six prefixes: 1) the ones with prefix $\langle a \rangle$, 2) the ones with prefix $\langle b \rangle$, .. and 3) the ones with prefix $\langle f \rangle$.

Find Subsets of Sequential Patterns

The subsets of sequential patterns can be mined by constructing the corresponding set of projected databases and mining each recursively. The projected databases as well as sequential patterns found in them are listed in Table 2, while the mining process is explained as follows:

- Find sequential patterns with prefix $\langle a \rangle$. Only the sequences containing $\langle a \rangle$ should be collected. Moreover, in a sequence containing $\langle a \rangle$, only the subsequence prefixed with the first occurrence of $\langle a \rangle$ should be considered. For example, in sequence $\langle (ef)(ab)(df)(cb) \rangle$, only the subsequence $\langle (_b)(df)(cb) \rangle$ should be considered for mining sequential patterns prefixed with $\langle a \rangle$. Notice that $(_b)$ means that the last element in the prefix, which is a, together with b, form one element.

The sequences in S containing $\langle a \rangle$ are projected with regards to $\langle a \rangle$ to form the $\langle a \rangle$ -projected database, which consists of four suffix sequences: $\langle (abc)(ac)d(cf) \rangle$, $\langle (d)c(bc)(ae) \rangle$, $\langle (b)(df)cb \rangle$, and $\langle (_f)cbc \rangle$.

By scanning the $\langle a \rangle$ -projected database once, its locally frequent items are $a : 2$, $b : 4$, $b : 2$, $c : 4$, $d : 2$, and $f : 2$. Thus, all the length-2 sequential patterns prefixed with $\langle a \rangle$ are found, and they are: $\langle aa \rangle : 2$, $\langle ab \rangle : 4$, $\langle (ab) \rangle : 2$, $\langle ac \rangle : 4$, $\langle ad \rangle : 2$, and $\langle af \rangle : 2$.

Recursively, all sequential patterns with prefix $\langle a \rangle$ can be partitioned into six subsets:

1) Those prefixed with $\langle aa \rangle$, 2) those with $\langle ab \rangle$, .., and, finally, 3) those with $\langle af \rangle$. These subsets can be mined by constructing respective projected databases and mining each recursively as follows:

- The $\langle aa \rangle$ -projected database consists of two nonempty (suffix) subsequences prefixed with $\langle aa \rangle$: $\langle (bc)(ac)d(cf) \rangle$, $\langle (_e) \rangle$. Since there is no hope to generate any frequent subsequence from this projected database, the processing of the $\langle aa \rangle$ -projected database terminates.
- The $\langle ab \rangle$ -projected database consists of three suffix sequences: $\langle (_c) \rangle(ac)d(cf) \rangle$, $\langle (_c)a \rangle$, and $\langle c \rangle$. Recursively mining the $\langle ab \rangle$ -projected database returns four sequential patterns:
 - $\langle (_c) \rangle$, $\langle (_c)a \rangle$, $\langle a \rangle$, and $\langle c \rangle$ (i.e., $\langle a(bc) \rangle$, $\langle a(bc)a \rangle$, $\langle aba \rangle$, and $\langle abc \rangle$.) They form the complete set of sequential patterns prefixed with $\langle ab \rangle$.
- The $\langle (ab) \rangle$ -projected database contains only two sequences: $\langle (_c)(ac)d(cf) \rangle$ and $\langle (df)cb \rangle$, which leads to the finding of the following sequential patterns prefixed with $\langle (ab) \rangle$: $\langle c \rangle$, $\langle d \rangle$, $\langle f \rangle$, and $\langle dc \rangle$.
- The $\langle ac \rangle$, $\langle ad \rangle$, and $\langle af \rangle$ -projected databases can be constructed and recursively mined similarly. The sequential patterns found are shown in Table 2.
 - Find sequential patterns with prefix $\langle b \rangle$, $\langle c \rangle$, $\langle d \rangle$, $\langle e \rangle$, and $\langle f \rangle$, respectively. This can be done by constructing the $\langle b \rangle$, $\langle c \rangle$, $\langle d \rangle$, $\langle e \rangle$, and $\langle f \rangle$ -projected databases and mining them, respectively. The projected databases as well as the sequential patterns found are shown in Table 2.

The Set Of Sequential Patterns is the Collection of Patterns Found in the above Recursive Mining Process.

One can verify that it returns exactly the same set of sequential patterns as what GSP and FreeSpan do.

Based on the above discussion, the algorithm of Prefix-Span is presented as follows:

Algorithm (PrefixSpan) Prefix-projected Sequential Pattern mining.

Input: A sequence database S , and the minimum support threshold min_support .

Output: The complete set of sequential patterns.

Method: Call $\text{PrefixSpan}(\langle \rangle, 0, S)$

Subroutine $\text{PrefixSpan}(\alpha, l, S|_{\alpha})$

The parameters are 1) α is a sequential pattern; 2) l is the length of α ; and 3) $S|_{\alpha}$ is the α -projected database if $\alpha \neq \langle \rangle$ otherwise, it is the sequence database S .

METHODS

- Scan $S|_{\alpha}$ once, find each frequent item, b , such that
 - b can be assembled to the last element of α to form a sequential pattern; or
 - $\langle b \rangle$ can be appended to α to form a sequential pattern.
- 2. For each frequent item b , append it to α to form a sequential pattern α' , and output α' .
- 3. For each α' , construct α' -projected database $S|_{\alpha'}$, and call $\text{PrefixSpan}(\alpha', l + 1, S|_{\alpha'})$.

High level design gives an overview of the logical flow. However this suffices the user to understand the logic.

Figure 1 depicts the basic knowledge about the system design and the architecture.

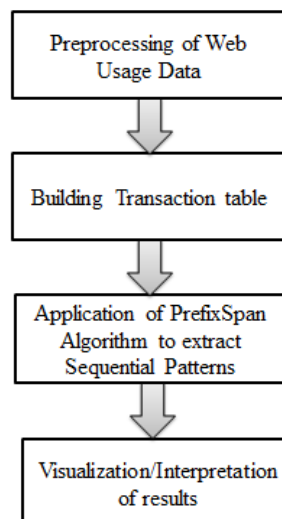


Figure 1: High Level Flow Chart

System Architecture of the proposed work is shown in figure 2. It consists of Preprocessing, Sequential Pattern Mining and Visualization of results modules.

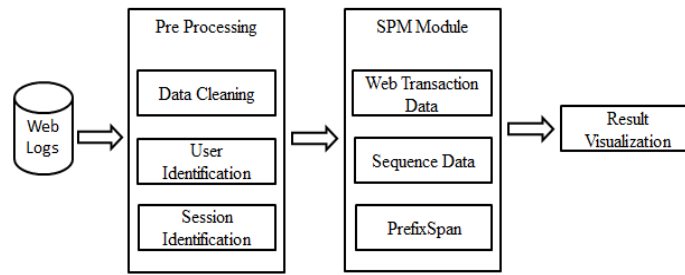


Figure 2: System Architecture

We analyze the efficiency of the algorithm as follows:

- No candidate sequence needs to be generated by PrefixSpan. Unlike a priori-like algorithms, PrefixSpan only grows longer sequential patterns from the shorter frequent ones. It neither generates nor tests any candidate sequence nonexistent in a projected database.
- Projected databases keep shrinking. A projected database is smaller than the original one because only the suffix subsequences of a frequent prefix are projected into a projected database. The shrinking factors can be significant because 1) usually, only a small set of sequential patterns grow quite long in a sequence database and, thus, the number of sequences in a projected database usually reduces substantially when prefix grows; and 2) projection only takes the whole string (not just suffix) and, thus, the shrinking factor is less than that of PrefixSpan.
- The major cost of PrefixSpan is the construction of projected databases. In the worst case, PrefixSpan constructs a projected database for every sequential pattern. If there exist a good number of sequential patterns, the cost is nontrivial. Techniques for reducing the number of projected databases will be discussed in the next subsection.

IV. EXPERIMENTAL RESULTS

To evaluate the effectiveness and efficiency of the PrefixSpan algorithm, we performed performance study of the algorithm by varying min_sup on data sets, with various kinds of sizes.

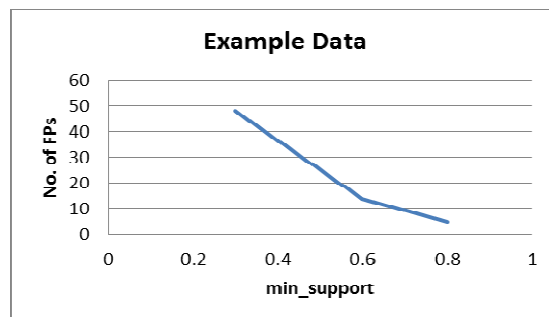


Figure 3: Distribution of Frequent Patterns of Example Data Set

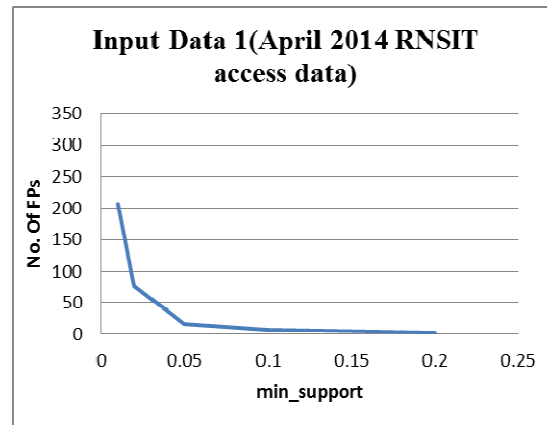


Figure 4: Distribution of Frequent Patterns for April 2014 RNSIT Access Data

CONCLUSIONS

We have performed a systematic study on mining of sequential patterns in large databases and developed a pattern-growth approach for efficient and scalable mining of sequential patterns.

Instead of refinement of the a priori-like, candidate generation-and-test approach, such as GSP, we promote a divide-and-conquer approach, called pattern-growth approach, which is an extension of FP-growth, an efficient pattern-growth algorithm for mining frequent patterns without candidate generation. Sequential pattern mining has broad applications including web reorganization, personalized marketing, customer retention.

ACKNOWLEDGEMENTS

The Author would like to thank the management of R N S Institute of Technology- HOD Dept. CSE, Dr G T Raju, Principal Dr. M K Venkatesha, and Director Dr. H N Shivashankar for the encouragement.

REFERENCES

1. R. Agrawal and R. Srikant, "Fast Algorithms for Mining Association Rules," Proc. 1994 Int'l Conf. Very Large Data Bases(VLDB '94), pp. 487-499, Sept. 1994.
2. R. Agrawal and R. Srikant, "Mining Sequential Patterns," Proc.1995 Int'l Conf. Data Eng. (ICDE '95), pp. 3-14, Mar. 1995.
3. R.J. Bayardo, "Efficiently Mining Long Patterns from Databases,"Proc. 1998 ACM-SIGMOD Int'l Conf. Management of Data (SIGMOD '98), pp. 85-93, June 1998.
4. R.J. Bayardo, R. Agrawal, and D. Gunopulos, "Constraint-Based Rule Mining on Large, Dense Data Sets," Proc. 1999 Int'l Conf. DataEng. (ICDE '99), pp. 188-197, Apr. 1999.
5. C. Bettini, X.S. Wang, and S. Jajodia, "Mining Temporal Relationships with Multiple Granularities in Time Sequences," Data Eng.Bull., vol. 21, pp. 32-38, 1998.
6. S. Guha, R. Rastogi, and K. Shim, "Rock: A Robust Clustering Algorithm for Categorical Attributes,"

- Proc. 1999 Int'l Conf. DataEng. (ICDE '99), pp. 512-521, Mar. 1999.
7. J. Han, G. Dong, and Y. Yin, "Efficient Mining of Partial Periodic Patterns in Time Series Database," Proc. 1999 Int'l Conf. Data Eng.(ICDE '99), pp. 106-115, Apr. 1999.
 8. J. Han, J. Pei, B. Mortazavi-Asl, Q. Chen, U. Dayal, and M.-C. Hsu, "FreeSpan: Frequent Pattern-Projected Sequential Pattern Mining," Proc. 2000 ACM SIGKDD Int'l Conf. Knowledge Discovery in Databases (KDD '00), pp. 355-359, Aug. 2000.

