An Inverted-ITL Algorithm for Mining Partial Periodic-Frequent Patterns

Ye-In Chang

Dept. of Computer Science and Engineering National Sun Yat-Sen University Kaohsiung, Taiwan, Republic of China changyi@mail.cse.nsysu.edu.tw Xin-Long Chen

Dept. of Computer Science and Engineering National Sun Yat-Sen University Kaohsiung, Taiwan, Republic of China chenxl@db.cse.nsysu.edu.tw Sheng-Hsin Chiang Dept. of Computer Science and Engineering National Sun Yat-Sen University Kaohsiung, Taiwan, Republic of China chiangsh@db.cse.nsysu.edu.tw

Abstract—In this paper, we propose an Inverted-ITL algorithm for mining partial periodic-frequent patterns. Although the GPFgrowth algorithm has been proposed to achieve the same goal, it needs to construct the prefix tree which needs too much storage space and processing time. Moreover, the GPF-growth algorithm will scan the database twice and sort each transactional list. To avoid those disadvantages, in this paper, we propose the inverted-ITL algorithm. Our algorithm is more efficient than the GPFgrowth algorithm. We will use one data structure, the ITL-tree, to store the items which appear in the database. Therefore, our algorithm can need shorter processing time than the GPF-growth algorithm.

Keywords—Frequent Patterns, Itemsets Mining, Periodic Patterns, Periodic-Frequent Patterns, Transactional Database

I. INTRODUCTION

In recent years, the research of data mining [1, 2] has become more popular. There are many techniques for mining interesting patterns in the database, like frequent pattern mining [3, 4, 5], frequent weighted pattern mining [6], frequent closed pattern mining [7], maximal frequent pattern mining [8], and periodicfrequent pattern mining [9, 10, 11, 12, 13]. Among the above techniques, current researches on periodic-frequent pattern mining have focused on discovering full periodic-frequent patterns in the database. However, partial periodic-frequent patterns are more common in the real world, *i.e.*, finding frequent patterns which frequently occur but do not successively occur one after one in the database. The reason for this phenomenon is the imperfect nature of the real-world.

The periodic-frequent pattern mining is an extension of the frequent pattern mining. In the frequent pattern mining, the count of the support is the factor which we only care about. However, in the periodic-frequent pattern mining, we do not just care about the frequency of each item, but also need to make sure that the frequent pattern occurs periodically in the transactional database [9]. According to the property of the transactional database, the database consists of two fields: (1) timestamp, and (2) item. The timestamp denotes the current time, when a transaction list is generated. For the topic of mining periodic-frequent patterns, the item will not appear repeatedly in a transaction list. Moreover, we do not have to care about the number of the items in a transaction list. That is, a transaction list {*a*, *a*, *b*, *c*} will not appear. Moreover, {*a*, *b*, *c*} and {*b*, *c*, *a*} are the same. The timestamp is used to record the time, when a

transaction list is generated and the transaction lists are sorted according to the ascending order of the timestamp. Therefore, for the entire database, there may exist time that no transaction list is generated. For instance, the list of transactions may contain [T1, T2, T3, T5, T6], where time T4 does not appear.

Take Table 1 as an example, which is denoted by transactional database TDB1. Assume that the user-specified minimum support is 2, and the user-specified maximum period is 3. At first, we scan all the transactions of transactional database TDB1 to count the total support and obtain the maximum period of the complete set of periods for each item. We can exploit these data items to discover which patterns are periodic-frequent patterns. Fig. 1 shows the result after scanning database TDB1, which is denoted by database TDB2. In database TDB2, the symbols 'ItemN', 'CountN' and 'MaxD' denote the item-name, the count of the total support, and the maximum period of the complete set of periods for each item, respectively.

TABLE I. AN EXAMPLE OF THE TRANSACTIONAL DATABASE TDE	TABLE I.	AN EXAMPLE OF THE TRANSACTIONAL DATABASE TE	ЭB	1
---	----------	---	----	---

ts	Items	ts	Items
1	ab	6	degf
2	acdi	7	abi
3	cefij	8	cdej
4	abfgh	9	abef
5	bd	10	acgi

According to the data of transactional database TDB2, we can know the support and the maximum period of each item in the database. Take item 'a' as an example. CountN is larger than or equal to 2 and MaxD is lower than or equal to 3. So item 'a' is a periodic-frequent pattern. Let's focus on the item 'c', which CountN is not less than 2 but its MaxD is larger than 3, so item 'c' is not a periodic-frequent pattern. Next, for item 'i', it is not a periodic-frequent pattern for the same reason as item 'c'. Although its CountN is larger than 2, its MaxD is larger than 3. Then, items 'j' and 'g' are not periodic-frequent patterns, because they have the same reason as the items 'c' and 'i'. That is, the MaxD is larger than the threshold. Finally, item 'h' also is not a periodic-frequent pattern. Because its CountN and MaxD, neither of them meet the threshold. Based on the definition of *anti-monotone*, pruning these items will not affect

the result. Therefore, the result after the pruning step contains items 'a', 'b', 'd', 'e ', and 'f'.

ItemN	CountN	MaxD	! : larger than 3
√a	6	3	!! : lower than 2
b	5	3	* : not periodic-frequent patt
*c	4	5!	
d	4	3	
*i	4	4!	
e	4	3	
f	4	3	
*j	2	5!	
*g	3	4!]
*h	1!!	6!	

Fig. 1. After scanning all transactions of database TDB1, called database TDB2 (the minimum support = 2, the maximum period = 3)

Next, we combine those items into the new candidate size 2 patterns. Then, after checking, the qualified patterns are 'ab' and 'ef'. In the same way, we combine the periodic-frequent size 2 patterns into the new candidate size 3 patterns to discover the periodic-frequent size 3 patterns. The similar process is repeated until no periodic-frequent patterns are generated.

Recently, Kiran *et al.* have proposed the PFP-growth++ algorithm [10], a prefix tree based algorithm. They use the prefix tree with the unique timestamp to mine the periodic-frequent patterns. However, they do not take the real world situation into account. Because in their algorithm, as long as the pattern does not meet the threshold once, it will not be considered as the periodic-frequent pattern. Later, Kiran *et al.* proposed the GPFgrowth algorithm [12], which is also a prefix tree based algorithm. The GPF-growth algorithm [12] considers fault tolerance by adding parameter Periodic-Ratio. It can effectively solve the problems encountered by the PFP-growth++ algorithm [10].

The GPF-growth algorithm [12] can find more realistic patterns than the PFP-growth++ algorithm [10]. The GPFgrowth algorithm considers about the fault tolerance of each pattern. However, we think that the construction of the prefix tree of GPF-growth needs too much time and memory space. Because the GPF-growth algorithm will scan database twice, and reorder each transactional list. Moreover, it constructs multiple prefix trees. Therefore, in this paper, we propose the Inverted-ITL algorithm to reduce the processing time. In our algorithm, we just only scan database once and store information of each pattern in two data structures. Later, we use these structures to find the periodic-frequent patterns. Therefore, our algorithm can need less processing time than the GPF-growth algorithm. Note that we do not need too much time to sort each transactional list and construct multiple prefix trees. From our performance study, we show that the performance of our algorithm is more efficient than that of the GPF-growth algorithm.

II. RELATED WORKS

The Apriori algorithm [3] has a great start and contribution to the research of data mining with frequent patterns. However, the algorithm is hard to achieve good performance. Later, Han et al. propose the FP-growth algorithm [5] for mining frequent patterns which scans twice in the database, and the algorithm exploits a tree structure, called FP-tree with the count of each item in the FP-tree. Next, Deng et al. propose the PrePost algorithm [4] and the dFIN algorithm [4] for mining frequent patterns with the property of scanning the database twice. Moreover, the algorithm exploits a tree structure, called PPCtree which is an extension of the FP-tree and a vertical data structure, called N-list for each item in the PPC-tree. There are many algorithms for data mining which are developed based on this structure, such as the PFP-growth algorithm [14], the NAFCP algorithm [7], the INLA-MFP algorithm [8], and the NFWI algorithm [6]. Although many algorithms have been proposed as mentioned above, none of them have considered the period. Kiran et al. have proposed the PFP-growth++ algorithm [10] for mining periodic-frequent patterns, which considers both the frequency and the period for each item. Later, Kiran and Reddy use the simplified model [11] to find all frequent patterns which have exhibited complete cyclic repetitions in the database. Kiran et al. propose the MCPF-model [13] to discover periodicfrequent patterns involving both frequent and rare items effectively. Kiran et al. [9, 10] have discussed greedy search techniques to discover periodic-frequent patterns effectively. All of these researches have focused on finding full periodicfrequent patterns. Then, Kiran et al. propose the GPF-growth algorithm [12], which considers about the fault tolerance for each item. The whole algorithm is basically the same as the PFPgrowth++ algorithm. The difference is the way to decide whether a pattern is a periodic-frequent pattern, which will be determined according to the percentage of the interest period. The components that make up the GPF-growth algorithm are the GPF-list and the GPF-tree. The GPF-growth algorithm forms the GPF-list by scanning the database once. Then, the GPF-list is sorted in the descending order of support for each item. Finally, they recursively mine the GPF-tree to discover the complete set of partial periodic-frequent patterns. Therefore, the GPF-growth algorithm needs to scan the database twice, sort the items once and the database once, and build a tree structure to find complete set of the partial periodic-frequent patterns.

III. THE INVERTED ITL ALGORITHM

In this section, we will introduce our algorithm called the Inverted-ITL (Inverted-Item-Time-List) algorithm to mine the partial periodic-frequent patterns.

ts	Items	ts	Items
1	ac	8	cdfg
2	abg	9	ab
3	de	10	cdef
4	abcd	11	abcef
6	abcd	12	abcd
7	abe	13	cef

 TABLE II.
 AN EXAMPLE OF THE TRANSACTIONAL DATABASE TDB3

A. Data Structure

In this subsection, we will use an example database TDB3 to explain our algorithm. During the process of partial periodicfrequent pattern mining, we exploit the support, the interesting period, and the periodic-ratio as main factors to determine whether the pattern is the partial periodic-frequent pattern or not. Thus, we use the *ITL-tree*, and the *ITL-list* to record the information for each item. Note that we only have to scan the database once to record the information and we do not need to do any sorting operation with the database. During the mining process, we will frequently use these data structures to find the partial periodic-frequent patterns. Table 2 shows an example database for timestamp *ts*, and Table 3 shows the variables used in our algorithm. Next, we will illustrate the two data structures which we will use them in our algorithm.

TABLE III. VARIABLES

Variable	Definition							
ItemN	The name of the pattern							
CountN	The Count of the pattern							
RangeC	The number of the pattern which is not larger than MaxPer during the period							
FirstT	The First Timestamp of the pattern							
LastT	The Last Timestamp of the pattern							
BP	Representing the timestamp of the bit pattern							
Mark	A boolean flag to check whether the last period of the pattern is not larger than the MaxPer							
MinSup	The Minimum Support threshold							
MinPR	The Minimum Periodic-Ratio threshold							
MaxPer	The Maximum Period threshold							

The *ITL-tree*, a tree structure, is used to store all the patterns which appear in the database. It is a prefix tree, where each node may become a partial periodic-frequent pattern, but the root is an empty pattern. All *k*-length patterns are stored at level *k* of the tree, where $k \ge 1$. For instance, in Fig. 2, square boxes represent non-candidate patterns, dotted circles represent the candidate patterns, and solid circles represent the partial periodic-frequent patterns. Moreover, in the *ITL-tree*, once the pattern is confirmed as a non-candidate pattern, its super-set pattern will not appear. In other words, no pattern will generate new patterns with non-candidate patterns. Finally, the patterns in the solid circle are the partial periodic-frequent patterns that we want to find.



Fig. 2. The simple diagram of ITL-tree

The *ITL-list* as shown in Fig. 3-(a) (which is the result after processing the second transaction) contains *ItemN*, *CountN*, *RangeC*, *FirstT*, *LastT*, *Mark*, and *BP* for each pattern. It is used to record the useful information of each pattern. Take Table 2 as an example. Assume that we have the *MinSup* = 3, the *MinPR* = 75%, and the *MaxPer* = 2. First, we scan on the first transaction,

"1: ac " with ts = 1. Since pattern 'a' and pattern 'c' are the first occurrence, so we create their own *ITL-list*. Because the ts (timestamp) of the first occurrence is not larger than the *MaxPer* and the *RangeC* plus 1. Therefore, we modify their *CountN*, *RangeC*, *FirstT*, and *LastT* to 1, 1, 1, 1, respectively. Since *FirstT* only records the ts of the first occurrence, so even if it appears later, we do not need to deal with this variable. Then, we have to deal with the *BP* part as shown in Fig. 3-(b), where *BP* is a table composed of a series of items and boolean values. For the above example, we know that pattern 'a' and pattern 'c' appear, when ts is 1. So in their *BP* table, we will set the boolean value to *True* at T1.

Next, we scan the second transaction, "2: abg" with ts = 2. Since pattern 'b' and pattern 'g' are the first occurrence, so we create their own *ITL-list*. Because ts of the first occurrence is equal to the *MaxPer* and the *RangeC* plus 1. Therefore, we modify their *CountN*, *RangeC*, *FirstT* and *LastT* to 1, 1, 2, 2, respectively. Then, the *CountN*, *RangeC*, and *LastT* values of pattern 'a' are updated to 2, 2, and 2, respectively. Because the current ts minus the *LastT* of pattern 'a' is less than *MaxPer*, *RangeC* is increased by 1. At the same time, for pattern 'a', we set T2 in *BP* to *True*, and for patterns 'b' and 'g', we do the same thing. The result is shown in Fig. 3-(a) and Fig. 3-(b).

					-		
Item N	a	c	b	g	Item N	T1	T2
CountN	2	1	1	1	a	Т	Т
RangeC	2	1	1	1	С	Т	
FirstT	1	1	2	2	b		Т
LastT	2	1	2	2	g		Т
	. (a)			. ((b)	

Fig. 3. *ITL-list* and the *BP* table after scanning the second transaction ('*abg*') : (a) the *ITL-list*; (b) the *BP* table.

The same step is processed until the entire database is scanned, and the variables which have not been mentioned, *Mark*, will not be decided until the whole database has been scanned. For instance, through Table 2, we can know that pattern 'a' last appears in *ts* 12. If the difference between the timestamp of the last transaction of the database and *LastT* of the pattern is less than or equal to *MaxPer*, we set *Mark* to *True* and increase *RangeC* by 1. Take pattern 'a' as an example. The gap between 13 (*ts* of the last transaction) and 12 (*LastT* of 'a') is less than *MaxPer*, so *Mark* of 'a' is set to *True* and we increase its *RangeC* by 1. The complete *ITL-list* of each pattern at level 1 is shown in Fig. 4.

B. The Mining Process

In this subsection, we will illustrate how to construct the *ITLtree* and the *ITL*-*list* in details. Moreover, we will illustrate the mining process of our algorithm.

According to the given transactional database TDB3, which is shown in Table 2, we let MinSup = 3, MinPR = 75%, and MaxPer = 2 as thresholds to illustrate our algorithm. First of all, we scan database TDB3 once to construct the *ITL-tree* and the *ITL-list*. *ITL-list* is a list which belongs exclusively to each node on the *ITL-tree*. Therefore, in our algorithm, we mainly use *ITL*-

tree to mine the partial periodic-frequent pattern and the *ITL-list* is used to record the useful information of each pattern.

ItemN	CountN	RangeC	FirstT	LastT 12	
a	8	9	1		
c	8	8	1	13	
b	7	8	2	12 8 12 13	
g	2	1	2		
d	6	6	3		
e	5	3	3		
f	4	4	8	13	

ItemN	T1	T2	T3	T4	T5	T6	T7	T8	T9	T10	T11	T12	T13	Mark
а	Т	Т		Т		Т	Т		Т		Т	Т		Т
с	Т			Т		Т		Т		Т	Т	Т	Т	Т
b		Т		Т		Т	Т		Т		Т	Т		Т
g		Т						Т						F
d			Т	Т		Т		Т		Т		Т		Т
e			Т				Т			Т	Т		Т	Т
f								Т		Т	Т		Т	Т

Fig. 4. The complete ITL-list of each pattern at level 1

Let's start with a transactional database TDB3 to construct the ITL-tree as shown in Fig. 5 (which is the result after processing the second transaction) and the ITL-list. Basically, there are two cases which we must concern: (1) the item never appears before; (2) the item has appeared before. For the first case, the item never appears before, we add a new node at level 1 of the *ITL-tree* and give it a name. At the same time for this node, we create its ITL-list. For the second case, the item has appeared before, we need to find the node with this name from level 1 in the ITL-tree, and then update its ITL-list. When we scan the first transaction, we can know that there are two patterns in the transaction, which are pattern 'a' and 'c'. When the algorithm reads the items in a transaction, it will change the ITLlist of these patterns in the ITL-tree. For instance, pattern 'a' is the first pattern to be scanned, so we insert node 'a' into level 1 of the ITL-tree. At the same time, we will create an ITL-list belonging to pattern 'a' and start to modify the data inside. The *CountN* is set to 1, because pattern 'a' appears for the first time. The *RangeC* is set to 1, because its period (*i.e.*, 1 - 0 = 1, where 0 is the starting time) is less than *MaxPer*. The *FirstT* is set to 1, because the first occurrence is at ts 1. The LastT is set to 1, because the last occurrence is at ts 1. Next, we set the T1 of BP table to True. The remaining items in the first transaction is executed by the same way.

Next, we scan the second transaction. The second transaction has patterns 'a', 'b', and 'g'. Since pattern 'a' already exists at level 1 of the *ITL-tree*, we only need to update the data of its *ITL-list*. Take pattern 'a' as an example. In the *ITL-list* of pattern 'a', its *CountN* is set to 2. Because it has appeared once before, we have 1 + 1 = 2. The *RangeC* is set to 2, because its period (*i.e.*, 2 - 1 = 1, where the first 1 is LastT) is less than *MaxPer*. The *LastT* is set to 2, because the pattern 'a' appears at *ts* 2 now. Then, for T2 of *BP* table, we set its *BP* table entry to True. Next, we turn to 'b' and 'g'. Take pattern 'b' as an example. Since pattern 'b' appears for the first time, we insert node 'b' into level 1 of the *ITL-tree* and create its *ITL-list*. Then, we set its *CountN* to 1. The *RangeC* is set to 1, because its period (*i.e.*, 2 - 0 = 2, where 0 is the starting time) is equal to *MaxPer*. The *FirstT* is set to 2, because the first occurrence is at *ts* 2. The *LastT* is set to 2, because the last occurrence is at *ts* 2. Moreover, we set the T2 of *BP* table to True. Since 'b' does not appear at *ts* 1, T1 of 'b' is empty. The condition of pattern 'g' is the same as pattern 'b', so we do what has done for pattern 'b' once for pattern 'g'. As a result, the *ITL-tree* and the *ITL-list* after scanning the second transaction is shown in Fig. 5.



Fig. 5. The ITL-tree and the ITL-list after scanning the second transaction

Finally, we perform the above process until the entire database has been scanned to complete level 1 of the *ITL-tree*, and the *ITL-list* are shown in Fig. 6. But we have not gotten the last period from the last *ts* (13). This last period will determine whether we have to increase *RangeC* by 1 and set *Mark* to *True*. At the same time, we can also know which of the patterns at level 1 are non-candidate patterns, candidate patterns and partial periodic-frequent patterns.



ItemN	CountN	RangeC	FirstT	LastT	BP
а	8	8	1	12	
ItemN	CountN	RangeC	FirstT	LastT	BP
c	8	7	1	13	
ItemN	CountN	RangeC	FirstT	LastT	BP
b	7	7	2	12	
ItemN	CountN	RangeC	FirstT	LastT	BP
g	2	1	2	8	
ItemN	CountN	RangeC	FirstT	LastT	BP
d	6	5	3	12	
ItemN	CountN	RangeC	FirstT	LastT	BP
e	5	2	3	13	
ItemN	CountN	RangeC	FirstT	LastT	BP
f	4	3	8	13	

Fig. 6. The ITL-tree and the ITL-list after scanning the entire database

Therefore, the modified *ITL-tree* and the modified *ITL-list* are shown in Fig. 7. Through Fig. 7, according to the *ITL-list* of each node, we can know, whether this node can become the partial periodic-frequent pattern. As long as the number of *RangeC* in the *ITL-list* is not less than $MinPR \times (MinSup+1)$ and *CountN* is not less than MinSup, we set it to the dotted circle. Furthermore, if *RangeC* is not less than $MinPR \times (CountN+1)$, we set it to the solid circle, and this pattern is the partial periodic-frequent pattern which we need. The remaining square boxes are considered unnecessary, because its *CountN* is less than *MinSup* or its *RangeC* is less than *MinPR* $\times (MinSup+1)$. So even if it is merged with other patterns, it cannot be a partial periodic-frequent size 1 patterns are patterns 'a', 'c', 'b', 'd', and 'f'.

Since we have found all partial periodic-frequent size 1 patterns, we can start finding partial periodic-frequent patterns of size 2. First, we will exploit these dotted or solid circle patterns in Fig. 7, where *RangeC* is larger than or equal to the *MinPR* × (*MinSup*+1) and *CountN* is not less than *MinSup*. We merge two patterns on the same level of the



ItemN	CountN	RangeC	FirstT	LastT	BP	Mark
a	8	9	1	12		Т
ItemN	CountN	RangeC	FirstT	LastT	BP	Mark
с	8	8	1	13	222	Т
ItemN	CountN	RangeC	FirstT	LastT	BP	Mark
b	7	8	2	12		Т
ItemN	CountN	RangeC	FirstT	LastT	BP	Mark
g	2	1	2	8	- 574	F
ItemN	CountN	RangeC	FirstT	LastT	BP	Mark
d	6	6	3	12		Т
ItemN	CountN	RangeC	FirstT	LastT	BP	Mark
e	5	3	3	13		Т
ItemN	CountN	RangeC	FirstT	LastT	BP	Mark
f	4	4	8	13		Т

Fig. 7. The modified ITL-tree and the modified ITL-list

ITL-tree with the same prefix into a (level+1) super-set patterns. Note that the prefix can be null and we do not need to merge with the square box. Because it is impossible to generate the partial periodic-frequent pattern. Since we want to merge the two patterns, we need to decide the order of the merged patterns. In the process of generating the size 2 patterns, we will merge the patterns at level 1 from left to right. Take Fig. 7 as an example. Patterns at level 1 of the

ITL-tree have the order from left to right: ['a', 'c', 'b', 'd', 'e', 'f']. Moreover, the pattern will generate a new pattern with each pattern after its order, and then it will be the next pattern. For instance, pattern 'a' will generate the new size 2 pattern with pattern 'c', 'b', 'd', 'e', 'f' and then it will turn to pattern 'c' and pattern 'b', 'd', 'e', 'f' to generate the new size 2 pattern, and so on.

Next, we illustrate how to generate size 2 patterns at level 2 of the *ITL-tree* and complete their *ITL-list*. Take pattern 'a' and 'f' as an example. At first, since we cannot find the pattern 'af' connected to 'a', we create a node 'af' and connect it to 'a' with an edge. Then, we compare the *FirstT* of patterns 'a' and 'f' to find the larger value. This value means that when we check the BP table of two patterns, we only need to check from the boolean value at this position. Because it will never appear at the same timestamp before. Similarly, we compare the *LastT* of the patterns 'a' and 'f' to find the smaller value. This value means that when we check the BP table of two patterns, we only need to check the boolean value at this position from the front. Because after that, they will never appear at the same timestamp. Take Fig. 8 as an example. From this Fig., we can know that pattern 'fdoes not appear before ts 8 and 'a' does not appear after ts 12. Therefore, we only need to check the boolean value between ts 8 and ts 12, which avoids unnecessary actions.

						Begin ↓						End J	
	T1	T2	T3	T4	T5	T6	T7	T8	T9	T10	T11	T12	T13
a	Т	Т		Т		Т	Т		Т		Т	Т	
f								Т		Т	Т		Т

Fig. 8. The BP table for patterns 'a' and 'f'

Then, we do the 'AND' operation on BPs of patterns 'a' and 'f'. This BP table after the 'AND' operation is regarded as BP table of pattern 'af'. We will use this BP table to complete its ITL-list. At first, we need to find the timestamp with a *True* value between ts 8 and ts 12. Then, we know that the first ts with True is at ts 11. Therefore, we set its CountN, RangeC, FirstT, LastT to 1, 0, 11, and 11, respectively. Then, we find that the following timestamp has no True. Therefore, as long as we find the last period to modify the value of RangeC and Mark, we complete the ITL-list of 'af'. At the same time, we can confirm whether pattern 'af' is the noncandidate pattern, the candidate pattern or the partial periodic-frequent pattern through the ITL-list, which is shown in Fig. 9. We apply the above steps to the remaining patterns. The result is shown in Fig. 10. The partial periodicfrequent size 2 patterns are patterns 'ac', 'ab', 'cb', 'cd', 'ce', *'cf '*, and *'ef '*.



Fig. 9. The ITL-tree and the ITL-list after the pattern 'af' is inserted



Fig. 10. The size 2 patterns of the ITL-tree

After mining all the partial periodic-frequent size 2 patterns, we start to find the partial periodic-frequent size 3 patterns. Similarly, we will merge the two patterns with the same prefix at level 2 of the ITL-tree except for the non-candidate patterns. So based on the above sentence, we will insert the new size 3 patterns 'acb', 'cbd', 'cbe', 'cbf', 'cde', 'cbf', 'cef' into level 3 of the ITL-tree, and follow the previous steps to complete their ITL-list.

Then, we confirm the type of each pattern according to their *ITL-list*. The result is shown in Fig. 11. According to this figure, we can know that the pattern '*acb*' and '*cef*' are partial periodic-frequent patterns. Moreover, there is no size 4 pattern which can be generated. So all the partial periodic-frequent patterns in this database have been found.

Finally, we compare the difference between our Inverted-ITL algorithm and the GPF-growth algorithm [12] for mining partial periodic-frequent pattern mining. We consider the original data in Table 2 as the input. When our algorithm constructs the data structures for the mining process, our algorithm only needs to scan database once and does not need the sorting step. However, the GPF-growth algorithm needs scanning database twice and sorts each transaction. Moreover, the GPF-growth algorithm needs to generate a prefix tree based on the reordered database, and it will generate many prefix trees during the mining process.

IV. PERFORMANCE

A. The Performance Model

We compare the processing time of the *Inverted-ITL* algorithm and the *GPF-growth* algorithm [12]. for mining partial periodic-frequent patterns in the real and synthetic datasets. We will consider using different values of the minimum support threshold, the maximum period threshold, and the minimum periodic-ratio threshold to execute different size datasets on the algorithm. For the real datasets, we use the Retail dataset (http://fimi.ua.ac.be/data/) for experiments. The details of the dataset Retail contains the transaction count = 88162, the item count = 16470 and the average item count per transaction = 10.30. For the synthetic dataset, T10.I4.D100K was generated



Fig. 11. The size 3 patterns of the ITL-tree

by using the generator from the IBM Quest Dataset Generator. The parameter T, I, and D represents the average item count per transaction, the average maximal size of frequent itemsets, and the number of transactions in the dataset, respectively.

B. Experiments Results

First, let's deal with the Retail dataset. In Fig. 12, we show the comparison of the processing time for the Retail dataset and the synthetic dataset under the change of the maximum period threshold. In this experiment, we set the minimum support threshold = 0.01% and the minimum periodic-ratio threshold = 0.01%. Through Fig. 12-(a), we observe that both our algorithm and the *GPF-growth* algorithm maintain a fairly stable curve. The reason is that all items in the Retail dataset, except for the few specific items, the rest of the items will basically not repeat within the threshold, so this situation will happen. Moreover, according to this figure, it shows that we provide better performance than the *GPF-growth* algorithm which scans the dataset twice and sorts many times. From Fig. 12-(b), we observe that the performance of our algorithm is also better than that of the *GPF-growth* algorithm.



Fig. 12. A comparison of the processing time under the change of the maximum period threshold: (a) for the Retail dataset; (b) for the synthetic dataset.

In Fig. 13, we show the comparison of the processing time for the Retail dataset and the synthetic dataset under the change of the minimum support threshold. In this experiment, we set the maximum period threshold = 10 and the minimum periodic-ratio threshold = 0.01%. Through Fig. 13-(a), we observe that as the value of minimum support threshold increases, the processing time of the two algorithms decreases. The reason is that only a small number of patterns can meet the threshold of support, when the minimum support threshold increases. When the number of the candidate patterns continues to decrease, the processing time also decreases. Through this figure, we observe that when the number of candidate patterns continues to decrease, the time curve of the *GPF-growth* algorithm will change greatly. However, our algorithm changes relatively small. This is because as long as the number of candidate patterns increases, the number of times that the *GPF-growth* algorithm needs to be sorted and the number of prefix trees generated will also increase. However, our algorithm only needs to process the *ITL-list* of each candidate pattern which are already in the *ITL-tree*. Therefore, in Fig. 13-(a), the processing time of our algorithm is faster than the *GPF-growth* algorithm. In Fig. 13-(b), the performance of our algorithm is also better than that of the *GPF-growth* algorithm.



Fig. 13. A comparison of the processing time under the change of the minimum support threshold: (a) for the Retail dataset; (b) for the synthetic dataset.

In Fig. 14, we show the comparison of the processing time for the Retail dataset and the synthetic dataset under the change of the minimum periodic-ratio threshold. In this experiment, we set the maximum period threshold = 10 and the minimum support threshold = 0.01%. Through Fig. 14-(a), we can observe that the minimum periodic-ratio threshold has small effect on the performance of two algorithms. So our algorithm has better performance than the *GPF-growth* algorithm. In Fig. 14-(b), the performance of our algorithm is also better than that of the *GPF-growth* algorithm.

In Fig. 15, we show the comparison of the processing time for the Retail dataset and the synthetic dataset under the change of the data size. In this experiment, we set the maximum period threshold = 10, the minimum support threshold = 0.01%, and the minimum periodic-ratio threshold = 0.01%. Through Fig. 15-(a), we can observe that as the size of the dataset increases, the processing time also increases. However, the growth rate of the *GPF-growth* algorithm is significantly larger than that of our algorithm. Because as the size of the dataset increases, the *GPFgrowth* algorithm needs to sort the items in more transactions, but our algorithm only needs to deal with the *ITL-list* of each candidate pattern in the *ITL-tree*. Therefore, the *GPF-growth* algorithm needs a lot of time to deal with sorting and generating prefix trees, and we only need to deal with each pattern in the *ITL-tree*. Through Fig. 15-(b), the performance of our algorithm is also better than that of the *GPF-growth* algorithm.



(b)

Fig. 14. A comparison of the processing time under the change of the minimum periodic-ratio: (a) for the Retail dataset; (b) for the synthetic dataset.



(b)

70000

Data Size

85000

100000

Fig. 15. A comparison of the processing time under the change of the change of the data size: (a) for the Retail dataset; (b) for the synthetic dataset.

55000

40000

V. CONCLUSION

In this paper, we have proposed an *Inverted-ITL* algorithm which can efficiently mine the partial periodic-frequent patterns. In the data mining, we have constructed the *ITL-tree* and the *ITL-list* which need less processing time to store information of each pattern than the *GPF-growth* algorithm for the same transactional database. From our simulation result, we have shown that our algorithm is more efficient than the *GPF-growth* algorithm.

ACKNOWLEDGMENT

This research was supported in part by the National Science Council of Republic of China under Grant No. MOST 110-2221-E-110-054.

REFERENCES

- M. K. Gupta and P. Chandra, "A Comprehensive Survey of Data Mining," International Journal of Information Technology, Vol. 12, pp. 1243–1257, February 2020.
- [2] K. Kaithwas and P. Borkar, "A Review on Different Data Mining Algorithms and Selection Methods," Proc. of the 2019 Int. Conf. on Intelligent Sustainable Systems, pp. 511–515, February 2019.
- [3] R. Agrawal and R. Agrawal and R. Srikant, "Fast Algorithms for Mining Association Rules in Large Databases," Proc. of the 20th Int. Conf. on Very Large Data Bases, pp. 487–499, 1994.
- [4] Z.-H. Deng, "DiffNodesets: An Efficient Structure for Fast Mining Frequent Itemsets," Applied Soft Computing, Vol. 41, pp. 214–223, April 2016.
- [5] J. Han, J. Pei, and Y. Yin, "Mining Frequent Patterns without Candidate Generation," Proc. of the 2000 ACM SIGMOD Int. Conf. on Management of Data, Vol. 29, No. 2, pp. 1–12, 2000.
- [6] H. Bui, B. Vo, H. Nguyen, T.-A. Nguyen-Hoang, and T.-P. Hong, "A Weighted N-List-Based Method for Mining Frequent Weighted Itemsets," Expert Systems with Applications, Vol. 96, pp. 388–405, April 2018.
- [7] T. Le and B. Vo, "An N-List-Based Algorithm for Mining Frequent Closed Patterns," Expert Systems with Applications, Vol. 42, No. 19, pp. 6648–6657, Nov. 2015.
- [8] B. Vo, S. Pham, T. Le, and Z.-H. Deng, "A Novel Approach for Mining Maximal Frequent Patterns," Expert Systems with Applications, Vol. 73, pp. 178–186, May 2017.
- [9] R. U. Kiran and M. Kitsuregawa, "Novel Techniques to Reduce Search Space in Periodic-Frequent Pattern Mining," Proc. of the 19th Int. Conf. on Database Systems for Advanced Applications, pp. 377–391, April 2014.
- [10] R. U. Kiran, M. Kitsuregawa, and P. K. Reddy, "Efficient Discovery of Periodic-Frequent Patterns in Very Large Database," The Journal of Systems and Software, Vol. 112, pp. 110–121, Feb. 2016.
- [11] R. U. Kiran and P. K. Reddy, "Towards Efficient Mining of Periodic-Frequent Patterns in Transactional Databases," Proc. of the 21th Int. Conf. on Database and Expert Systems Applications, pp. 194–208, Aug./Sept. 2010.
- [12] R. U. Kiran, J. N. Venkatesh, M. Toyoda, M. Kitsuregawa, and P. K. Reddy, "Discovering Partial Periodic-Frequent Patterns in a Transactional Database," The Journal of Systems and Software, Vol. 125, pp. 170–182, March 2017.
- [13] A. Surana, R. U. Kiran, and P. K. Reddy, "An Efficient Approach to Mine Periodic- Frequent Patterns in Transactional Databases," Proc. of the 15th Pacific-Asia Conf. on Knowledge Discovery and Data Mining, pp. 254– 266, May 2011.
- [14] S.-S. Chen, T. C.-K. Huang, and Z.-M. Lin, "New and Efficient Knowledge Discovery of Partial Periodic Patterns with Multiple Minimum Supports," The Journal of Systems and Software, Vol. 84, No. 10, pp. 1638–1651, Oct. 2011.

Authors' background

Your Name	Title*	Research Field	Personal website
Ye-In Chang	Full professor	Data Mining	https://db.cse.nsysu.edu.tw/~changyi/
Xin-Long Chen	Master student	Data Mining	https://db.cse.nsysu.edu.tw/~chenxl/
Sheng-Hsin Chiang	Master student	Data Mining	https://db.cse.nsysu.edu.tw/~chiangsh/

*This form helps us to understand your paper better, the form itself will not be published.

*Title can be chosen from: master student, Phd candidate, assistant professor, lecture, senior lecture, associate professor, full professor