

A Space Optimization for FP-Growth

Eray Özkural and Cevdet Aykanat
Department of Computer Engineering
Bilkent University 06800 Ankara, Turkey
{erayo,aykanat}@cs.bilkent.edu.tr

Abstract

Frequency mining problem comprises the core of several data mining algorithms. Among frequent pattern discovery algorithms, FP-GROWTH employs a unique search strategy using compact structures resulting in a high performance algorithm that requires only two database passes. We introduce an enhanced version of this algorithm called FP-GROWTH-TINY which can mine larger databases due to a space optimization eliminating the need for intermediate conditional pattern bases. We present the algorithms required for directly constructing a conditional FP-Tree in detail. The experiments demonstrate that our implementation has a running time performance comparable to the original algorithm while reducing memory use up to twofold.

1. Introduction

Frequency mining is the discovery of all frequent patterns in a transaction or relational database. Frequent pattern discovery comprises the core of several data mining algorithms such as association rule mining and sequence mining [10], dominating the running time of these algorithms. The problem involves a transaction database $T = \{X | X \subseteq I\}$ that consists in a set of transactions each of which are drawn from a set I of items. The mining algorithm finds all patterns that occur with a frequency satisfying a given absolute support threshold ϵ . In practice, the number of items $|I|$ is in the order of magnitude of 10^3 and more. The number of transactions is much larger, at least 10^5 . A pattern is $X \subseteq I$, a subset of I , and the set of all patterns is 2^I . The frequency function $f(T, x) = |\{x \in Y | Y \in T\}|$

computes the number of times a given item $x \in I$ occurs in the transaction set T ; it is extended to sets of items $f(T, X) = |\{X \subseteq Y | Y \in T\}|$ to compute the frequency of a pattern.

Frequency mining is the discovery of all frequent patterns in a transaction set with a frequency of support threshold ϵ and more. The set of all frequent patterns is $\mathcal{F}(T, \epsilon) = \{X | X \subseteq 2^I \wedge f(T, X) \geq \epsilon\}$. In the algorithms, the set of frequent items $F = \{x \in I | f(T, x) \geq \epsilon\}$ may require special consideration. A significant property of frequency mining known as downward closure states that if $X \in \mathcal{F}(T, \epsilon)$ then $\forall Y \subset X, Y \in \mathcal{F}(T, \epsilon)$ [2].

An inherent limitation of frequency mining is the amount of main memory available [8]. In this paper, we present a space optimization to FP-Growth algorithm and we demonstrate its impact on performance with experiments on synthetic and real-world datasets. In the next section, we give the background on the FP-GROWTH algorithm. Section 3 and Section 4 explain our algorithm and implementation. Section 5 presents the experiments, following that we offer our conclusions.

2. Background

2.1. Compact structures

Compact data structures have been used for efficient storage and query/update of candidate item sets in frequency mining algorithms. SEAR [12], SPEAR [12], and DIC [6] use tries (also known as prefix trees) while FP-GROWTH [10] uses FP-Tree which is an enhanced trie structure.

Using concise structures can reduce both running time and memory size requirements of an algorithm. Tries are well known structures that are widely used

for storing strings and have decent query/update performance. The aforementioned algorithms exploit this property of the data structure for better performance. Tries are also efficient in storage. A large number of strings can be stored in this dictionary type which would not otherwise fit into main memory. For frequency mining algorithms both properties are critical as our goal is to achieve efficient and scalable algorithms. In particular, the scalability of these structures is quite high [10] as they allow an algorithm to track the frequency information of the candidate patterns for very large databases. The FP-Tree structure in FP-GROWTH allows the algorithm to maintain *all* frequency information in the main memory obtained from two database passes. Using the FP-Tree structure has also resulted in novel search strategies.

A notable work on compact structures is [15] in which a binary-trie based summary structure for representing transaction sets is proposed. The trie is further compressed using Patricia tries. Although significant savings in storage and improvements in query time are reported, the effectiveness of the scheme in a frequency mining algorithm remains to be seen. In another work in FIMI 2003 workshop [3], an algorithm called PATRICIAMINE using Patricia tries has been proposed [13]. The performance of PATRICIAMINE has been shown to be consistently good in the extensive benchmark studies of FIMI workshop [3]; it was one of the fastest algorithms although it was not the most efficient. For many applications, the average case performance may be more important than performing well in a small number of cases, therefore further research on this PATRICIAMINE would be worthwhile.

In this paper, we introduce an optimized version of FP-GROWTH. A closer analysis of it is in order.

2.2. FP-Growth algorithm

The FP-GROWTH algorithm uses the frequent pattern tree (FP-Tree) structure. FP-Tree is an improved trie structure such that each itemset is stored as a string in the trie along with its frequency. At each node of the trie, *item*, *count* and *next* fields are stored. The *items* of the path from the root of the trie to a node constitute the item set stored at the node and the *count* is the frequency of this item set. The node link *next* is a pointer to the next node with the same *item* in the FP-Tree. Field *parent* holds a pointer to the parent node, *null* for root. Additionally, we maintain a header table

which stores heads of node links accessing the linked list that spans all same items. FP-Tree stores only frequent items. At the root of the trie is a *null* item, and strings are inserted in the trie by sorting item sets in a unique¹ decreasing frequency order [10].

Table 1 shows a sample transaction set and frequent items in descending frequency order. Figure 1 illustrates the FP-Tree of sample transaction set in Table 1. As shown in [10], FP-Tree carries complete information required for frequency mining and in a compact manner; the height of the tree is bounded by maximal number of frequent items in a transaction. MAKE-FP-TREE (Algorithm 1) constructs an FP-Tree from a given transaction set T and support threshold ϵ as described.

Transaction	Ordered Frequent Items
$t_1 = \{f, a, c, d, g, i, m, p\}$	$\{f, c, a, m, p\}$
$t_2 = \{a, b, c, f, l, m, o\}$	$\{f, c, a, b, m\}$
$t_3 = \{b, f, h, j, o\}$	$\{f, b\}$
$t_4 = \{b, c, k, s, p\}$	$\{c, b, p\}$
$t_5 = \{a, f, c, e, l, p, m, n\}$	$\{f, c, a, m, p\}$

Table 1. A sample Transaction Set

In Algorithm 3 we describe FP-GROWTH which has innovative features such as:

1. Novel search strategy
2. Effective use of a summary structure
3. Two database passes

FP-GROWTH turns the frequency k -length pattern mining problem into “a sequence of k -frequent 1-item set mining problems via a set of conditional pattern bases” [10]. It is proposed that with FP-GROWTH there is “no need to generate any combinations of candidate sets in the entire mining process”. With an FP-Tree *Tree* given as input the algorithm generates all frequent patterns. There are two points in the algorithm that should be explained: the single path case and conditional pattern bases. If an FP-Tree has only a single path, then an optimization is to consider all combinations of items in the path (single path case is the ba-

¹ All strings must be inserted in the same order; the order of items with the same frequency must be the same.

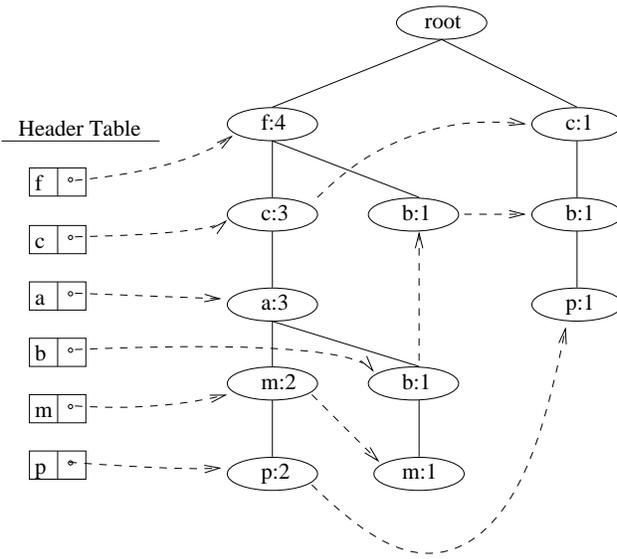


Figure 1. An FP-Tree Structure.

sis of recursion in FP-GROWTH). Otherwise, the algorithm constructs for each item a_i in the header table a *conditional pattern base* and an FP-Tree $Tree_\beta$ based on this structure for recursive frequency mining. Conditional pattern base is simply a compact representation of a derivative database in which only a_i and its prefix paths in the original $Tree$ occur. Consider path $\langle f : 4, c : 3, a : 3, m : 2, p : 2 \rangle$ in $Tree$. For mining patterns including m in this path, we need to consider only the prefix path of m since the nodes after m will be mined elsewhere (in this case only p). In the prefix path $\langle f : 4, c : 3, a : 3 \rangle$ any pattern including m can have frequency equal to the frequency of m , therefore we may adjust the frequencies in the prefix path as $\langle f : 2, c : 2, a : 2 \rangle$ which is called a *transformed prefix path* [10]. The set of transformed prefix paths of a_i forms a small database of patterns which co-occur with a_i and thus contains complete information required for mining patterns including a_i . Therefore, recursively mining conditional pattern bases for all a_i in $Tree$ is equivalent to mining $Tree$ (which is equivalent to mining the complete DB.). $Tree_\beta$ in FP-GROWTH is the FP-Tree of the conditional pattern base.

FP-GROWTH is indeed remarkable with its unique divide and conquer approach. Nevertheless, it may be profitable for us to view it as generating candidates despite the title of [10]. The conditional pattern base is

a set of candidates among which only some of them turn out to be frequent. The main innovation however remains intact: FP-GROWTH takes advantage of a tailored data structure to solve the frequency mining problem with a divide-and-conquer method and with demonstrated efficiency and scalability. Besides, the conditional pattern base is guaranteed to be smaller than the original tree, which is a desirable property. An important distinction of this algorithm is that, when examined within the taxonomy of algorithms, it employs a unique search strategy. When the item sets tested are considered, it is seen that this algorithm is neither DFS nor BFS. The classification for FP-GROWTH in Figure 3 of [11] may be slightly misleading. As Hipp later mentions in [11], “FP-Growth does not follow the nodes of the tree . . . , but directly descends to *some part* of the itemsets in the search space”. In fact, the part is so well defined that it would be unfair to classify FP-GROWTH as conducting a DFS. It does not even start with item sets of small length and proceed to longer item sets. Rather, it considers a *set of patterns at the same time* by taking advantage of the data structure. This unique search strategy makes it hard to classify FP-GROWTH in the context of traditional uninformed search algorithms.

Algorithm 1 MAKE-FP-TREE(DB, ϵ)

- 1: Compute F and $f(x)$ where $x \in F$
 - 2: Sort F in frequency decreasing order as L
 - 3: Create root of an FP-Tree T with label “null”
 - 4: **for all** transaction $t_i \in T$ **do**
 - 5: Sort frequent items in t_i according to L . Let sorted list be $[p|P]$ where p is the head of the list and P the rest.
 - 6: INSERT-TRIE($[p|P]$)
 - 7: **end for**
-

3. An improved version of FP-Growth

During experiments with large databases, we have observed that FP-GROWTH was costly in terms of memory use. Thus, we have experimented with improvements to the original algorithm. In this section, we propose FP-GROWTH-TINY (Algorithm 4) which is an enhancement of FP-GROWTH featuring a space optimization and minor improvements. An important optimization eliminates the need for intermediate conditional pattern bases. A minor improvement comes

Algorithm 2 INSERT-TRIE($[p|P], T$)

```
1: if  $T$  has a child  $N$  such that  $item[N] = item[p]$ 
   then
2:    $count[N] \leftarrow count[N] + 1$ 
3: else
4:   Create new node  $N$  with  $count = 1$ ,  $parent$ 
     linked to  $T$  and node-link linked to nodes with
     the same item via  $next$ 
5: end if
6: if  $P \neq \emptyset$  then
7:   INSERT-TRIE( $P, N$ )
8: end if
```

Algorithm 3 FP-GROWTH($Tree, \alpha$)

```
1: if  $Tree$  contains a single path  $P$  then
2:   for all combination  $\beta$  of the nodes in path  $P$  do
3:     generate pattern  $\beta \cup \alpha$  with support minimum
       support of nodes in  $\beta$ 
4:   end for
5: else
6:   for all  $a_i$  in header table of  $Tree$  do
7:     generate pattern  $\beta \leftarrow a_i \cup \alpha$  with  $support =$ 
        $support[a_i]$ 
8:     construct  $\beta$ 's conditional pattern base and
       then  $\beta$ 's conditional FP-Tree  $Tree_\beta$ 
9:     if  $Tree_\beta \neq \emptyset$  then
10:      FP-GROWTH( $Tree_\beta, \beta$ )
11:    end if
12:   end for
13: end if
```

from not outputting all combinations of the single path in the basis of recursion. Instead, we output a representation of this task since subsequent algorithms can take advantage of a compact representation for generating association rules and so forth.² Another improvement is pruning the infrequent nodes of the single path.

In the following subsection, the space optimization is discussed.

3.1. Eliminating conditional pattern base construction

The conditional tree $Tree_\beta$ can be constructed directly from $Tree$ without an intermediate conditional pattern base. The conditional pattern base in

² For the FIMI workshop, we output all patterns separately as required. It can be argued that a meaningful mining of all frequent patterns must output them one by one.

Algorithm 4 FP-GROWTH-TINY($Tree, \alpha$)

```
1: if  $Tree$  contains a single path  $P$  then
2:   prune infrequent nodes of  $P$ 
3:   if  $|P| > 0$  then
4:     output "all patterns in  $2^P$  and  $\alpha$ "
5:   end if
6: else
7:   for all  $a_i$  in header table of  $Tree$  do
8:      $Tree_\beta \leftarrow$  CONS-CONDITIONAL-FP-TREE( $Tree, a_i$ )
9:     output pattern  $\beta \leftarrow a_i \cup \alpha$  with  $count =$ 
        $f(a_i) \triangleright f(x)$  of  $Tree$ 
10:    if  $Tree_\beta \neq \emptyset$  then
11:      FP-GROWTH( $Tree_\beta, \beta$ )
12:    end if
13:   end for
14: end if
```

FP-GROWTH can be implemented as a set of *patterns*. A pattern in FP-GROWTH consists of a set of *symbols* and an associated *count*. With a counting algorithm and retrieval/insertion of patterns directly into the FP-Tree structure, we can eliminate the need for such a pattern base. Algorithm 5 constructs a conditional FP-Tree from a given $Tree$ and a symbol s for which the transformed prefix paths are computed.

The improved procedure first counts the symbols in the conditional tree without generating an intermediate structure and constructs the set of frequent items. Then, each transformed prefix path is computed as patterns retrieved from $Tree$ and are inserted in $Tree_\beta$.

COUNT-PREFIX-PATH presented in Algorithm 6 scans the prefix paths of a given node. Since the pattern corresponding to the transformed prefix path has the count of the node, it simply adds the count to the count of all symbols in the prefix path. This step is required for construction of a conditional FP-Tree directly since an FP-Tree is based on the decreasing frequency order of F . This small algorithm allows us to compute the counts of the symbols in the conditional tree in an efficient way, and was the key observation in making the optimization possible.

Algorithm 7 retrieves a transformed prefix path for a given node excluding node itself and Algorithm 8 inserts a pattern into the FP-Tree. GET-PATTERN computes the transformed prefix path as described in [10]. INSERT-PATTERN prunes the items not present in the frequent item set F of $Tree$ (which does not have to be identical to the F of calling procedure) and sorts the pattern in decreasing frequency order to maintain FP-

Algorithm 5 CONS-CONDITIONAL-FP-TREE($Tree, s$)

```
1:  $table \leftarrow itemtable[Tree]$ 
2:  $list \leftarrow table[symbol]$ 
3:  $Tree' \leftarrow MAKE-FP-TREE$ 
4:  $\triangleright$  Count symbols without generating an intermediate structure
5:  $node \leftarrow list$ 
6: while  $node \neq null$  do
7:   COUNT-PREFIX-PATH( $node, count[Tree]$ )
8:    $node \leftarrow next[node]$ 
9: end while
10: for all  $sym \in range[count]$  do
11:   if  $count[sym] \geq \epsilon$  then
12:      $F[Tree'] \leftarrow F[Tree'] \cup sym$ 
13:   end if
14: end for
15:  $\triangleright$  Insert conditional patterns to  $Tree_\beta$ 
16:  $node \leftarrow list$ 
17: while  $node \neq null$  do
18:    $pattern \leftarrow GET-PATTERN(node)$ 
19:   INSERT-PATTERN( $Tree', pattern$ )
20:    $node \leftarrow next[node]$ 
21: end while
22: return  $Tree'$ 
```

Algorithm 6 COUNT-PREFIX-PATH($node, count$)

```
1:  $prefixcount \leftarrow count[node]$ 
2:  $node \leftarrow parent[node]$ 
3: while  $parent[node] \neq null$  do
4:    $count[symbol[node]] \leftarrow$ 
      $count[symbol[node]] + prefixcount$ 
5:    $node \leftarrow parent[node]$ 
6: end while
```

Algorithm 7 GET-PATTERN($node$)

```
1:  $pattern \leftarrow MAKE-PATTERN$ 
2: if  $parent[node] \neq null$  then
3:    $count[pattern] \leftarrow count[node]$ 
4:    $currnode \leftarrow parent[node]$ 
5:   while  $parent[currnode] \neq null$  do
6:      $symbols[pattern] \leftarrow symbols[pattern] \cup$ 
        $symbol[currnode]$ 
7:      $currnode \leftarrow parent[currnode]$ 
8:   end while
9: else
10:   $count[pattern] \leftarrow 0$ 
11: end if
12: return  $pattern$ 
```

Tree properties and adds the obtained string to the FP-Tree structure. The addition is similar to insertion of a single string, with the difference that insertion of a pattern is equivalent to insertion of the symbol string of the pattern $count[pattern]$ times.

Algorithm 8 INSERT-PATTERN($Tree, pattern$)

```
1:  $pattern \leftarrow pattern \cap F[Tree]$ 
2: Sort pattern in a predetermined frequency decreasing order
3: Add the pattern to the structure
```

The optimization in Algorithm 5 makes FP-GROWTH more efficient and scalable by avoiding additional iterations and cutting down storage requirements. An implementation that uses an intermediate conditional pattern base structure will scan the tree once, constructing a linked list with transformed prefix paths in it. Then, it will construct the frequent item set from the linked list, and in a second iteration insert all transformed prefix paths with a procedure similar to INSERT-PATTERN. Such an implementation would have to copy the transformed prefix paths twice, and iterate over all prefix paths three times, once in the tree, and twice in the conditional pattern list. In contrast, our optimized procedure does not execute any expensive copying operations and it needs to scan the pattern bases only twice in the tree. Besides efficiency, the elimination of extra storage requirement is significant because it allows FP-GROWTH to mine more complicated data sets with the same amount of memory.

An idea similar to our algorithm was independently explored in FP-GROWTH* by making use of information in 2-items [9]. In their implementations, Grahne and Zhu have used strategies based on 2-items to improve running time and memory usage, and they have reported favorable performance, which has also been demonstrated in the benchmarks of the FIMI '03 workshop [3].

4. Implementation notes

We have made a straightforward implementation of FP-GROWTH-TINY and licensed it under GNU GPL for public use. It has been written in C++, using GNU g++ compiler version 3.2.2.

For variable length arrays, we used `vector<T>` in standard library. For storing transactions, patterns and other structures representable as strings we have used efficient variable length arrays. We used `set<T>` to store item sets in certain places where it would be fast to do so, otherwise we have used sorted arrays to implement sets.

No low level memory or I/O optimizations were employed.

A shortcoming of the pattern growth approach is that it does not seem to be very memory efficient. We store many fields per node and the algorithm consumes a lot of memory in practice.

The algorithm has a detail which required a special code: sorting the frequent items in a transaction according to an order L , in line 2 of Algorithm 1 and line 2 of Algorithm 8. For preserving FP-Tree properties all transactions must be inserted in the very same order.³ The items are sorted first in order of decreasing frequency and secondarily in order of indices to achieve a unique frequency decreasing order. Using this procedure, we are not obliged to maintain an L .

5. Performance study

In this section we report on our experiments demonstrating the performance of FP-GROWTH-TINY. We have measured the performance of Algorithm 3 and Algorithm 4 on a 2.4Ghz Pentium 4 Linux system with 1GB memory and a common 7200 RPM IDE hard disk. Both algorithms were run on four synthetic and five real-world databases with varying support threshold. The implementation of the original FP-GROWTH algorithm is due to Bart Goethals.⁴

We describe the data sets used for our experiments in the next two subsections. Following that, we present our performance experiments and interpret the results briefly, comparing the performance of the improved algorithm with the original one.

5.1. Synthetic data

We have used the association rule generator described in [2] for synthetic data. Synthetic databases in our evaluation have been selected from [17] and [16].

³ For patterns also in our implementation.

⁴ Goethals has made his implementation publicly available at <http://www.cs.helsinki.fi/u/goethals/>

These databases have been derived from previous studies [1, 2, 14]. Table 2 explains the symbols we use for denoting the parameters of association rule generator tool. The experimental databases are depicted in Table 3. In all synthetic databases, $|I|$ is 1000, and $|\mathcal{F}_{max}|$ is 2000. The original algorithm could not process T20.I6.D1137 in memory therefore the number of transactions was decreased to 450K.

$ T $	Number of transactions in transaction set
$ t _{avg}$	Average size of a transaction t_i
$ f_m _{avg}$	Average length of maximal pattern f_m
I	Number of items in transaction set
$ \mathcal{F}_{max} $	Number of maximal frequent patterns

Table 2. Dataset parameters

Name	$ T $	$ t _{avg}$	$ f_m _{avg}$
T10.I6.1600K	1.6×10^6	10	6
T10.I4.1024K	1.024×10^6	10	4
T15.I4.367K	3.67×10^5	15	4
T20.I6.450K	4.5×10^5	20	6

Table 3. Synthetic data sets

5.2. Real-world data

We have used five publicly available datasets in the FIMI repository. `accidents` is a traffic accident data [7]. `retail` is market basket data from an anonymous Belgian retail store [5]. The `bms-webview1`, `bms-webview2` and `bms-pos` datasets are from a benchmark study described in [4]. Some statistics of the datasets are presented in Table 4.

5.3. Memory consumption and running time

The memory consumption and running time of FP-GROWTH-TINY and FP-GROWTH are plotted for varying relative supports from 0.25% to 0.75% in Figure 2 and Figure 3 for synthetic databases and Figure 4 and Figure 5 for real-world databases except for accidents database which is a denser database that should be mined at 10% and more. The implementations were run

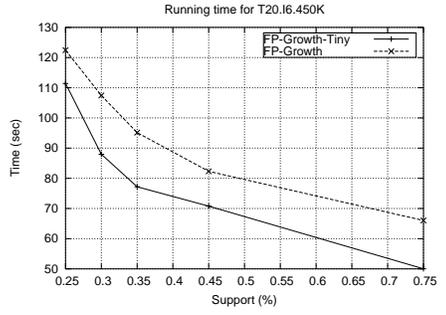
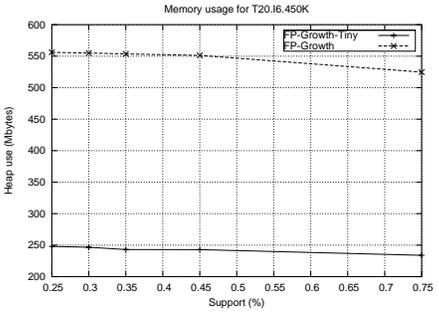
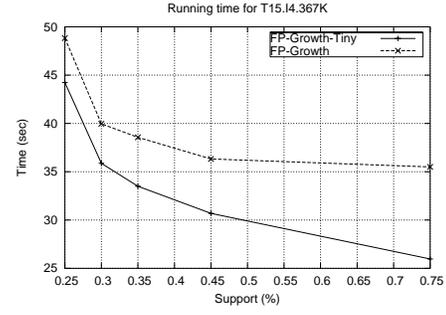
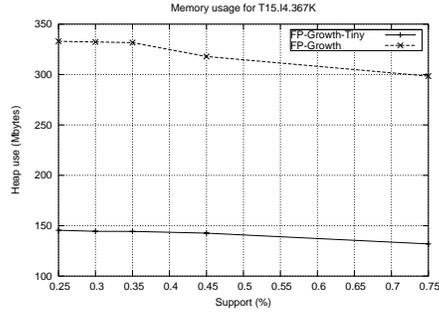
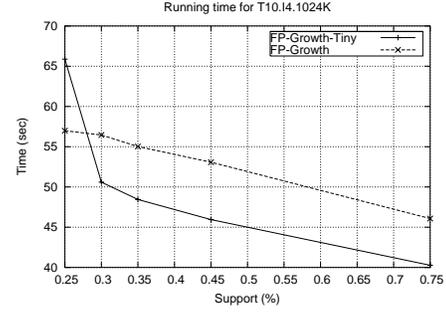
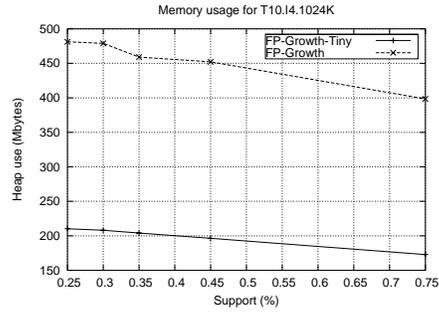
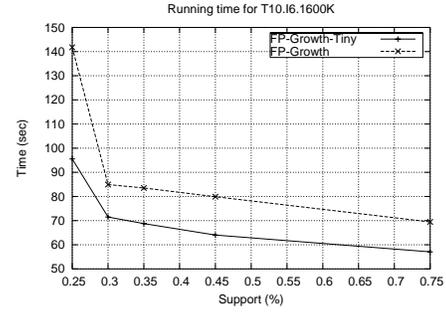
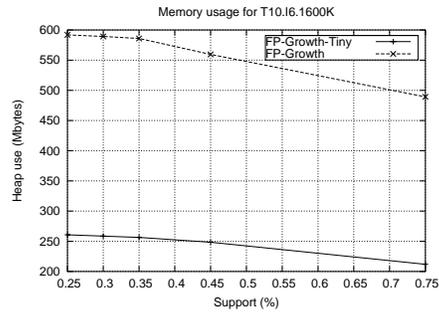


Figure 2. Memory consumption of FP-GROWTH-TINY and FP-GROWTH on synthetic databases

Figure 3. Running time performance of FP-GROWTH-TINY and FP-GROWTH on synthetic databases

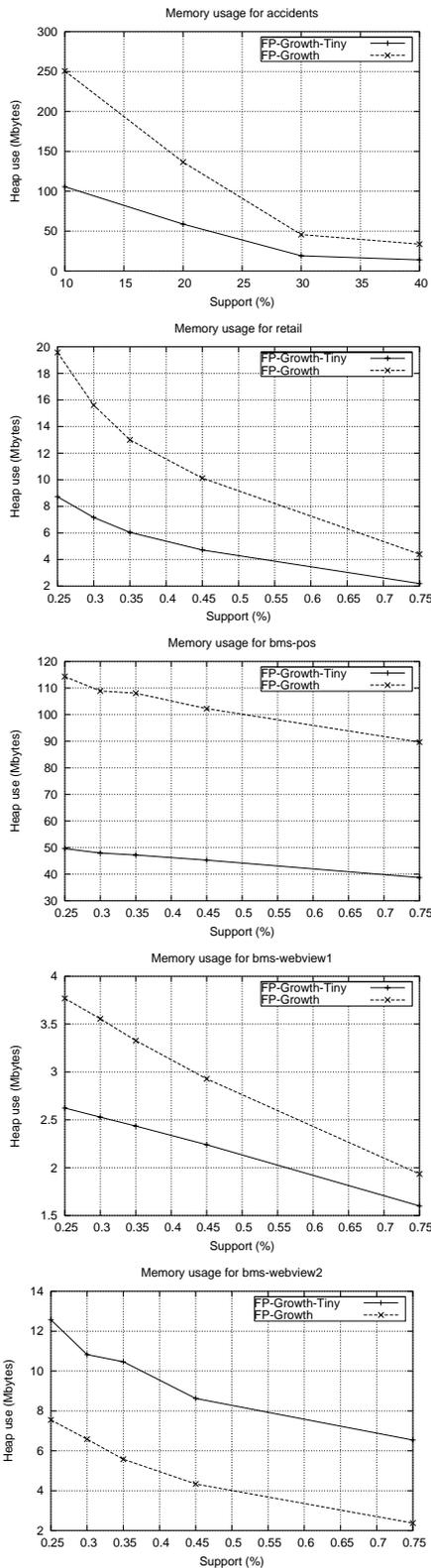


Figure 4. Memory consumption of FP-GROWTH-TINY and FP-GROWTH on real-world databases

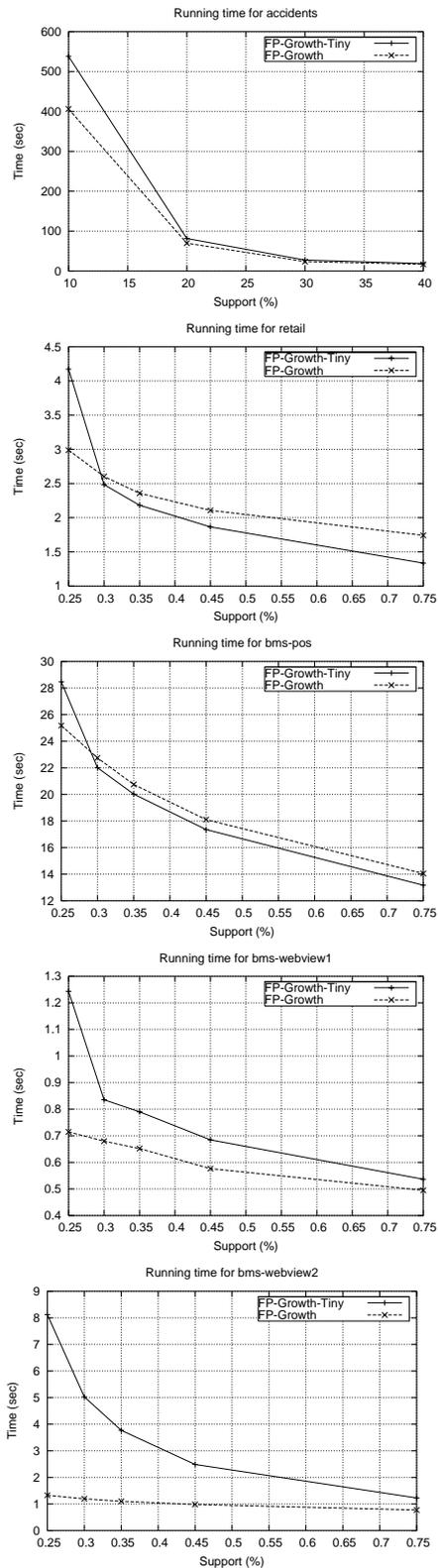


Figure 5. Running time performance of FP-GROWTH-TINY and FP-GROWTH on real-world databases

Name	$ T $	$ I $	$ t _{avg}$
accidents	3.41×10^5	469	33.81
retail	8.82×10^4	16470	10.31
bms-pos	5.16×10^5	1657	6.53
bms-webview1	5.96×10^4	60978	2.51
bms-webview2	7.75×10^4	330286	4.62

Table 4. Real-world data sets

inside a typical KDE desktop session. The running time is measured as the wall-clock time of the system call. The memory usage is measured using the GNU glibc tool `memusage`, considering only the maximum heap size since stack use is much smaller than heap size.

The plots for synthetic datasets are similar among themselves, while we observe more variation in real-world datasets. Memory is saved in all databases, except in `bms-webview2`, which requires 2.74 times the memory used in FP-GROWTH; this has an adverse effect on running time as discussed below. In others, we observe that memory usage reduces down to 41.5% in accidents database with 4% support, which is 2.4 times smaller than FP-GROWTH.

Due to the optimization, our implementation can process larger databases than the vanilla version. For most problem instances, the memory consumption has been reduced more than twofold compared to the original algorithm. An advantage of our approach is that with the same amount of memory, we can process more complicated databases.⁵ The experiments overall show that the conditional pattern base construction which we have eliminated has a significant space cost during the recursive construction of conditional FP-Trees.

The running time behaviors of two algorithms are quite similar on the average. Our algorithm tends to perform better and is faster in higher support thresholds, while in lower thresholds the performance gap becomes closer. FP-GROWTH-TINY runs faster except in `bms-webview1`, `bms-webview2` and lower thresholds of T10.I4.1024K. In `bms-webview1` database, FP-GROWTH-TINY runs 10-27% slower; in `bms-webview2` database we observe that FP-GROWTH-TINY has slowed down by a factor of 5.56 for 0.25% support threshold, and slowdown is observed also for other support thresholds (down to

⁵ Note that FP-GROWTH uses a compressed representation of frequency information, whose size may be thought of as related to complexity of the dataset.

50%). In T10.I4.1024K we see 12% slowdown for 0.25% support and 2% slowdown for 0.3% support. In other problem instances FP-GROWTH-TINY, runs faster, up to 28.5% for retail dataset at 0.75% support.

In the figures, we observe a relation between memory saving and decreased running time. We had expected that improving space utilization would remarkably decrease the running time. However, we have not observed as large an improvement as we would have liked in running time. On the other hand, our trials show significant improvement in memory use contrasted to vanilla FP-GROWTH, allowing us to mine more complicated/larger datasets with the same amount of memory.

The adverse situation with `bms-webview1` and `bms-webview2` shows that the performance study must be extended to determine whether the undesirable behavior recurs at a large scale, since these are both sparse data sets coming from the same source. At any rate, a closer inspection of FP-GROWTH-TINY seems necessary. We anticipate that the benchmark studies at the FIMI workshop will illustrate its performance more precisely.

6. Conclusions

We have presented our version of FP-GROWTH which sports multiple improvements in Section 3. An optimization over the original algorithm eliminates a large intermediate structure required in the recursive step of the published FP-GROWTH algorithm in addition to two other minor improvements.

In Section 5, we have reported the results of our performance experiments on synthetic and real-world databases. The performance of the optimized algorithm has been compared with a publicly available FP-GROWTH implementation. We have observed more than twofold improvement in memory utilization over the vanilla algorithm. In the best case, memory size has become 2.4 times smaller, while in the worst case memory saving was not possible in a small real-world database. Typically, our implementation makes better use of memory, enabling it to mine larger and more complicated databases that cannot be processed by the original algorithm. The running time behavior of both algorithms are quite similar on the average; FP-GROWTH-TINY runs up to 28.5% percent faster, however it may run slower in a minority of instances.

References

- [1] R. Agrawal and J. C. Shafer. Parallel mining of association rules. *IEEE Trans. On Knowledge And Data Engineering*, 8:962–969, 1996.
- [2] R. Agrawal and R. Srikant. Fast algorithms for mining association rules. In J. B. Bocca, M. Jarke, and C. Zaniolo, editors, *Proc. 20th Int. Conf. Very Large Data Bases, VLDB*, pages 487–499. Morgan Kaufmann, 12–15 1994.
- [3] M. J. Z. Bart Goethals. Fimi '03: Workshop on frequent itemset mining implementations. In *Proceedings of the ICDM 2003 Workshop on Frequent Itemset Mining Implementations*, Melbourne, Florida, USA, 2003.
- [4] Z. Z. Blue. Real world performance of association rule algorithms. In *KDD 2001*.
- [5] T. Brijs, G. Swinnen, K. Vanhoof, and G. Wets. Using association rules for product assortment decisions: A case study. In *Knowledge Discovery and Data Mining*, pages 254–260, 1999.
- [6] S. Brin, R. Motwani, J. D. Ullman, and S. Tsur. Dynamic itemset counting and implication rules for market basket data. In J. Peckham, editor, *SIGMOD 1997, Proceedings ACM SIGMOD International Conference on Management of Data, May 13-15, 1997, Tucson, Arizona, USA*, pages 255–264. ACM Press, 05 1997.
- [7] K. Geurts, G. Wets, T. Brijs, and K. Vanhoof. Profiling high frequency accident locations using association rules. In *Proceedings of the 82nd Annual Transportation Research Board*, page 18pp, Washington DC (USA), January 2003.
- [8] B. Goethals. Memory issues in frequent itemset mining. In *Proceedings of the 2004 ACM Symposium on Applied Computing (SAC'04)*, Nicosia, Cyprus, March 2004.
- [9] G. Grahne and J. Zhu. Efficiently using prefix-trees in mining frequent itemsets. In *Proceedings of the First IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'03)*.
- [10] J. Han, J. Pei, and Y. Yin. Mining frequent patterns without candidate generation. In W. Chen, J. Naughton, and P. A. Bernstein, editors, *2000 ACM SIGMOD Intl. Conference on Management of Data*, pages 1–12. ACM Press, May 2000.
- [11] J. Hipp, U. Güntzer, and G. Nakhaeizadeh. Algorithms for association rule mining – a general survey and comparison. *SIGKDD Explorations*, 2(1):58–64, July 2000.
- [12] A. Mueller. Fast sequential and parallel algorithms for association rule mining: A comparison. Technical Report CS-TR-3515, College Park, MD, 1995.
- [13] A. Pietracaprina and D. Zandolin. Mining frequent itemsets using patricia tries. In *Proceedings of the First IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'03)*.
- [14] A. Savasere, E. Omiecinski, and S. B. Navathe. An efficient algorithm for mining association rules in large databases. In *The VLDB Journal*, pages 432–444, 1995.
- [15] D.-Y. Yang, A. Johar, A. Grama, and W. Szpankowski. Summary structures for frequency queries on large transaction sets. In *Data Compression Conference*, pages 420–429, 2000.
- [16] M. J. Zaki, S. Parthasarathy, and W. Li. A localized algorithm for parallel association mining. In *ACM Symposium on Parallel Algorithms and Architectures*, pages 321–330, 1997.
- [17] M. J. Zaki, S. Parthasarathy, M. Ogihara, and W. Li. Parallel algorithms for discovery of association rules. *Data Mining and Knowledge Discovery*, 1(4):343–373, 1997.