

# Comparative Study of Frequent Itemset Mining Algorithms Apriori and FP Growth

Ritu Garg  
Student MTECH (CS)  
M.D.U Rohtak (Haryana, India)

Preeti Gulia, PhD  
Associate Prof.  
M.D.U Rohtak (Haryana, India)

## ABSTRACT

Frequent itemset mining leads to the discovery of associations among items in large transactional database. In this paper, two algorithms[7] of generating frequent itemsets are discussed: Apriori and FP-growth algorithm. In apriori algorithm candidates are generated and testing is done which is easy to implement but candidate generation and support counting is very expensive in this because database is checked many times. In the fp-growth, there is no candidate generation and requires only 2 passes over the database but in this the generation of fp-tree become very expansive to built and support is counted only when entire dataset is added to fp-tree. The comparison of these algorithms will tell which algorithm is better to perform.

## Keywords

Frequent itemset mining, Apriori, FP-Growth

## 1. INTRODUCTION

In recent years amount of data in the database has increased rapidly. The increasing size of the database has led to growing interest in extraction of useful information from the bulk of data. Data mining is a technique useful for attaining useful information from vast databases. Implicit information within a database can be very useful in tasks such as marketing, financial forecast etc. This information has to be derived efficiently. Frequent itemset mining discovers significant relationships among variables or items in a dataset.

Association rule mining[5] searches for relationships between items in a dataset. It finds association among set of items in transactional database. Each transaction is a list of items. Association rules[4] is in form  $A \Rightarrow B$  which means customer buys A also tends to buy B. To mine association rule, basic concepts of support and confidence are needed. **Support**  $s$  is the probability that a transaction contain (X, Y). **Confidence**  $C$  is the measure of the strength of the association rule, suppose the confidence of the association rule  $x \Rightarrow y$  is 90%, it means that 90% of the transactions that contain X also contain Y together. Also minimum support and minimum confidence is needed to eliminate the unimportant association rules. Such that the association rules is hold when it is greater than the minimum support and minimum confidence.

T_id	Items
100	a, b,c
200	a, c
300	a, d
400	b, e, f

Equation for support and confidence:  
Support ( $A \Rightarrow B$ ) =Probability ( $A \cap B$ ).  
Confidence ( $A \Rightarrow B$ ) =Probability ( $B/A$ ).

Let the min\_support and min\_confidence are 50%.for association rule  $a \Rightarrow c$ , support (a, c) = $2/4 * 100\% = 50\%$ . Confidence = Support (a, c)/Support (a) = $50\%/75\% = 66.6\%$ , means that customer buys a also have 66.6% chance to buy c.

## 2. APRIORI ALGORITHM

The apriori algorithm[2] is firstly proposed by R.Aggarwal and R.Srikant in 1994 for mining frequent itemset. In data mining, Apriori is a classic algorithm for learning association rules. Apriori is designed to operate on databases containing transactions (for example, collections of items bought by customers, or details of a website frequentation).

Apriori algorithm follows two phases:

- Generate Phase: In this phase candidate (k+1)-itemset is generated using k-itemset; this phase creates  $C_k$  candidate set.
- Prune Phase: In this phase candidate set is pruned to generate large frequent itemset using “minimum support” as the pruning parameter. This phase creates  $L_k$  large itemset

Fig 1 shows the pseudo code for apriori algorithm :

Apriori_Algo(L,C,k)	
<b>Pass 1</b>	
1.	Generate the candidate itemsets in $C_1$
2.	Save the frequent itemsets in $L_1$
<b>Pass k</b>	
1.	Generate the candidate itemsets in $C_k$ from the frequent itemsets in $L_{k-1}$
i.	Join $L_{k-1} p$ with $L_{k-1} q$ , as follows: <b>insert</b> <span style="float:right">into <math>C_k</math></span> <b>select</b> $p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$ <b>from</b> $L_{k-1} p, L_{k-1} q$ <b>where</b> $p.item_1 = q.item_1, \dots, p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$
ii.	Generate all (k-1)-subsets from the candidate itemsets in $C_k$
iii.	Prune all candidate itemsets from $C_k$ where some (k-1)-subset of the candidate itemset is not in the frequent itemset $L_{k-1}$
2.	Scan the transaction database to determine the support for each candidate itemset in $C_k$
3.	Save the frequent itemsets in $L_k$

**Example of apriori algorithm:**

- Consider a database, D , consisting of 9 transactions.
- Suppose min. support count required is 2 (i.e.  $\text{min\_sup} = 2/9 = 22\%$  )

TID	List of items
T100	I1 ,I2 ,I5
T200	I2 ,I4
T300	I2 ,I3
T400	I1 ,I2 ,I4
T500	I1 ,I3
T600	I2 ,I3
T700	I1 ,I3
T800	I1 ,I2 ,I3 ,I5
T900	I1 ,I2 ,I3

Table 2: Database consisting 9 transactions

**Step 1:** Count the number of transactions in which each item occurs (Table 3.a)

**Step 2:** In this step we remove all the items that are bought less than 2 times from the table (Table 3.b)

C1		L1		
Itemset	Sup-count		Itemset	Sup-count
I1	6	Compare and prune	I1	6
I2	7		I2	7
I3	6		I3	6
I4	2		I4	2
I5	2		I5	2

(A) (b)

Table 3: first scan of Apriori( Scan for count of each candidate)

**Step 3:** Make all the pairs of items by using property JOIN L1 with L1and count how many times each pair is bought together (Table 4.a)

**Step 4:** Remove all the item pairs with number of transactions less than two (Table 4.b)

(a)		Compare and prune	(b)	
Itemset	Sup-count		Itemset	Sup-count
I1,I2	4	Compare and prune	I1,I2	4
I1,I3	4		I1,I3	4
I1,I4	1		I1,I5	2
I1,I5	2		I2,I3	4
I2,I3	4		I2,I4	2
I2,I4	2		I2,I5	2
I2,I5	2		I3,I4	0
I3,I4	0		I3,I5	1
I3,I5	1		I4,I5	0
I4,I5	0			

(a) (b)

Table 4: The second scan of A-priori (Generate C2 and Scan D for count of each Candidate).

**Step 5:** To make the set of three items we need one more rule (it's termed as self-join),

It simply means, from the Item pairs in the above table, we find two pairs with the same first Item

C3		Compare And Prune	L3	
Itemset	Sup-count		Itemset	Sup-count
I1 ,I2 ,I3	2	Compare And Prune	I1 ,I2 ,I3	2
I1 ,I2 ,I5	2		I1 ,I2 ,I5	2

(a) (b)  
Table 5: The third scan of A-priori (Generate C3 and Scan D for count of each Candidate)

- While we are on this, suppose you have sets of 3 items say ABC, ABD, ACD, ACE, BCD and you want to generate item sets of 4 items you look for two sets having the same first two alphabets.

ABC and ABD -> ABCD  
ACD and ACE -> ACDE

**Step 6:** according to above statement I1, I2, I3, I5 is generated whose minimum support is less than 2.so this is not frequent.

Thus the set of three items that are bought together most frequently are I1, I2, I3 and I1, I2, I5

**ADVANTAGES:**

1. Use large itemset.
2. Easy to implement.
3. Easily parallelized.

**DISDVANTAGE:**

1. It may need to generate a huge no of candidate sets. So its generation is expensive.
2. Assumes transactional database is memory resident.
3. Support count is expensive because require many database scan.

**3. FP-GROWTH ALGORITHM**

The FP-Growth Algorithm[1], proposed by Han in, is an efficient and scalable method for mining the complete set of frequent patterns by pattern fragment growth, using an extended prefix-tree structure for storing compressed and crucial information about frequent patterns named frequent-pattern tree (FP-tree). In his study, Han proved that his method outperforms other popular methods for mining frequent patterns [1],[8],[9],[10], e.g. the Apriori Algorithm

Major steps in FP-growth is

Step1- It firstly compresses the database showing frequent item set in to FP-tree. FP-tree is built using 2 passes over the dataset.

Step2: It divides the FP-tree in to a set of conditional database and mines each database separately, thus extract frequent item sets from FP-tree directly. It consist of one root labeled as null, a set of item prefix sub trees as the children of the root, and a frequent .item header table. Each node in the item prefix sub tree consists of three fields: item-name, count and node link where--- item-name registers which item the node represents; count registers the number of transactions represented by the portion of path reaching this node, node link links to the next node in the FP- tree. Each item in the header table consists of two fields---item name and head of node link, which points to the first node in the FP-tree carrying the item name.

### 3.1 FP-Tree structure

The frequent-pattern tree (FP-tree)[6] is a compact structure that stores quantitative information about frequent patterns in a database. Han defines the FP-tree as the tree structure defined below:

1. One root labeled as “null” with a set of item-prefix subtrees as children, and a frequent-item-header table:
  - i. Each node in the item-prefix subtree consists of three fields: Item-name: registers which item is represented by the node;
  - ii. Count: the number of transactions represented by the portion of the path reaching the node;
  - iii. Node-link: links to the next node in the FP-tree carrying the same item-name, or null if there is none.

Each entry in the frequent-item-header table consists of two fields:

- i. Item-name: as the same to the node;
- ii. Head of node-link: a pointer to the first node in the FP-tree carrying the item-name.

Procedure of fp-growth algorithm:

Input: constructed FP-tree  
Output: complete set of frequent patterns  
Method: Call FP-growth (FP-tree, null).  
Procedure FP-growth (Tree,  $\alpha$ )  
{  

1. If Tree contains a single path P then
2. For each combination do generate pattern  $\beta \sqcup \alpha$  with support = minimum support of nodes in  $\beta$ .
3. Else for each header  $a_i$  in the header of Tree do {
4. Generate pattern  $\beta = a_i \sqcup \alpha$  with support =  $a_i$ .support;
5. Construct  $\beta$ .s conditional pattern base and then  $\beta$ .s conditional FP-tree Tree  $\beta$
6. If Tree  $\beta =$  null
7. Then call FP-growth (Tree  $\beta$ ,  $\beta$ );

}

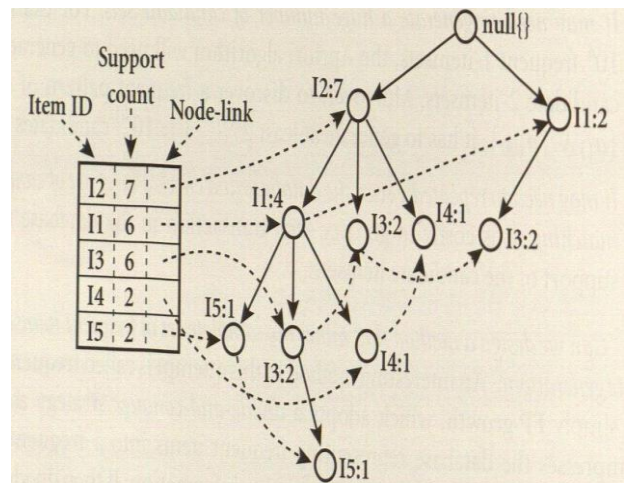
### Example:

Let us create the FP-tree for the example from Table 2:

- First we scan the database and determine the set of frequent items (1-itemsets) and their support counts(frequencies):  
 $L = \{ \{I2:7\}, \{I1:6\}, \{I3:6\}, \{I4:2\}, \{I5:2\} \}$
- Then we create the root of the FP-tree and label it with “null”
- We take each transaction, sort the items according to descending support count, and create a branch for it. For example the scan of the first transaction “T100:I1, I2, I5”, which contain tree items: I2, I1 and I5 in sorted descending, leads to the construction of the first branch of the tree: (I2:1), (I1:1), (I5:1).
- The second transaction T200 contains the items I2 and I4. This would result a branch where I2 is linked to the root and I4 is linked to I2. However this branch would share a common prefix, i2, with the existing path for T100. Therefore we instead increment the count of the I2 node by 1 and create a new node (I4:1), which is linked as a child of (I2:2).

In general when considering the branch to be added for a transaction, the count of each node along a common prefix is incremented by 1 and nodes for the items following the prefix are created and linked accordingly.

To facilitate tree traversal, an item header table is built so that each item points to its occurrences in the tree via a chain of node-links. In this way the problem of mining frequent pattern in database is transformed to that of mining the FP-tree.



The FP-tree is mined as follows: Start from each frequent length-1 pattern, as an initial suffix pattern, construct its conditional pattern base, a sub-database, which consists of the set of prefix paths in the FP-tree co-occurring with the suffix pattern, then construct its conditional FP-tree and perform mining recursively on such a tree. The pattern growth is achieved by the concatenation of the suffix pattern with the frequent patterns generated from a conditional FP-tree.

- Let us consider I5, which is the last item in L. I5 occurs in two branches of the FP-tree:
  - (I2, I1, I5:1)
  - (I2, I1, I3, I5:1)
- I5 is a suffix, so its corresponding two prefix paths are
  - (I2, I1:1)
  - (I2, I1, I3:1)
- Its conditional FP-tree contains only a single path: (I2:2, I1:2); I3 is removed because its support count of 1 is less than the minimum support count
- The single path generates all the combinations of frequent patterns:
  - {I2,I5:2}
  - {I1,I5:2}
  - {I2, I1, I5:2}
- For I4 exist 2 prefix path, which form the conditional pattern base:
  - {{I2, I1:1},{I2:1}}
- This generates a single-node conditional FP-tree:
  - (I2:2)
- The frequent pattern: {I2, I1:2}

The following table shows the frequent pattern generated for each node:

Item	Conditional Pattern Base	Conditional FP-tree	Frequent Pattern Generated
I5	{{I2, I1:1}, {I2, I1, I3:1}}	(I2:2, I1:2)	{I2, I5:2}, {I1, I5:2}, {I2, I1, I5:2}
I4	{{I2, I1:2}, {I2:1}}	(I2:2)	{I2, I4:2}
I3	{{I2, I1:2}, {I2:2}, {I1:2}}	(I2:4, I1:2), (I2:4)	{I2, I3:4}, {I1, I3:4}, {I2, I1, I3:2}, {I2, I1:4}
I1	{{I2:4}}	(I2:4)	{I2, I1:4}

**ADVANTAGES:**

1. It compresses the database.
2. Require only 2 pass over database.
3. There is no candidate generation.
4. Faster than apriori.
5. Reduces search cost

**DISADVANTAGE:**

1. It may not fit in main memory.
2. FP tree is expensive to build.
  - i. takes time to build but once built frequent itemset can be obtained easily.
  - ii. Support can only be calculated once the entire dataset is added to fp-tree.

#### 4. COMPARISON OF APRIORI AND FP-GROWTH ALGORITHMS

Parameters	Apriori Algorithm	FP-growth Algorithm
Technique	Use Apriori property and join and prune property	It constructs conditional frequent pattern tree and conditional pattern base from database which satisfy minimum support.
Memory utilization	Due to large no. of candidate generation require large memory space.	Due to compact structure and no candidate generation require less memory.
Number of scans	Multiple scans for generating candidate set.	Scan the database only twice and twice only.
Time	Execution time is more as time is wasted in producing candidate every time.	Execution time is lesser than the Apriori algorithm.

#### 5. CONCLUSION

Frequent itemset mining is an important task in association rule mining. It has been found useful in many applications like market basket analysis, financial forecasting etc. We have discussed about classical algorithm Apriori and Fp growth with their pros and cons. using the horizontal approach, owing to all candidate itemset for each level has to be discovered, the longer the length of the frequent itemset, more the number of candidate generation. Projected tree method is efficient in terms of speed but utilizes more space. These disadvantages can be overcome by using techniques like hashing, partitioning etc. In this paper study of itemset mining algorithms is done and on the basis of that study comparison is given between them.

#### 6. REFERENCES

- [1] J. Han, H. Pei, and Y. Yin. Mining Frequent Patterns without Candidate Generation. In: Proc. Conf. on the Management of Data (SIGMOD'00, Dallas, TX). ACM Press, New York, NY, USA 2000.
- [2] Agrawal, R. and Srikant, R. 1994. Fast algorithms for mining association rules. In Proc. 1994 Int. Conf. Very Large Data Bases (VLDB'94), Santiago, Chile, pp. 487–499.
- [3] Agarwal, R., Aggarwal, C., and Prasad, V.V.V. 2001. A tree projection algorithm for generation of frequent itemsets. Journal of Parallel and Distributed Computing, 61:350–371.
- [4] B.Santhosh Kumar and K.V.Rukmani. Implementation of Web Usage Mining Using APRIORI and FP Growth Algorithms. Int. J. of Advanced Networking and

Applications, Volume: 01, Issue: 06, Pages: 400-404 (2010).

- [5] Cornelia Gyorödi and Robert Gyorödi. A Comparative Study of Association Rules Mining Algorithms.
- [6] F. Bonchi and B. Goethals. FP-Bonsai: the Art of Growing and Pruning Small FP-trees. Proc. 8th Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD'04, Sydney, Australia), 155–160. Springer-Verlag, Heidelberg, Germany 2004.
- [7] Christian Borgelt. Keeping Things Simple: Finding Frequent Item Sets by Recursive Elimination. Workshop Open Source Data Mining Software (OSDM'05, Chicago, IL), 66-70. ACM Press, New York, NY, USA 2005
- [8] Aiman Moyaid, Said and P.D.D., Dominic and Azween, Abdullah. A Comparative Study of FP-growth Variations. International journal of computer science and network security, 9 (5). pp. 266-272.
- [9] Liu, G., Lu, H., Yu, J. X., Wang, W., & Xiao, X.. AFOPT: An Efficient Implementation of Pattern Growth Approach, In Proc. IEEE ICDM'03 Workshop FIMI'03, 2003.
- [10] Grahne, G., & Zhu, J. Fast Algorithm for frequent Itemset Mining Using FP-Trees. IEEE Transactions on Knowledge and Data Engineer, Vol.17, NO.10, 2005.