



Date :

Practical No. 8

Aim: To implement FP Growth algorithm.

Theory:

Frequent pattern Growth Algorithm:

This algo. is an improvement to the Apriori method. A FP is generated without the need for candidate generation. FP growth algo. represents the db in the form of tree called a FP tree or FP tree. This tree structure will maintain the association bet<sup>n</sup> the itemsets. The db is fragmented using one frequent item. This fragmented part is called 'pattern fragment'. The itemset of these fragmented pattern are analyzed. Thus with this method, the search for frequent item sets is reduced comparatively.

FP tree:

FP tree is a tree like structure that is made with the initial itemset of the database. The purpose of FP tree is to mine the most frequent pattern. Each node of FP tree represent an item of the interest itemset. The root node represents null while the lower nodes represent the itemsets. The association of the nodes with the lower nodes, that is the itemsets with the other itemsets are maintained while forming the tree.

Frequent pattern Algorithm steps:

- ① The first step is to scan the db to find the occurrences of the itemset in the database. This step is the same as the first step of Apriori. The count of 1-itemsets in the db is called support count or frequency of 1-itemset.
- ② The second step is to construct the FP tree. For this, create the root of the tree. The root is represented by null.



Date :

- ③ The next step is to scan the database again and examine the transaction. Examine the first transaction and findout the itemset in it. The itemset with the max count is taken at top, the next itemset with lower count and so on. It means that the branch of the tree is constructed with transaction item sets in descending order of count.
- ④ The next transaction in the database is examined. The item sets are ordered in descending order of count. If any itemset of this transaction is already present in another branch, then this transaction branch would share a common prefix to the root. This means that the common itemset is linked to the new node of another itemset in this transaction.
- ⑤ Also, the count of the itemset is incremented as it occurs in the transaction. Both the common node and new node count is increased by 1 as they are created and linked all to transaction.
- ⑥ The next step is to mine the created FP tree. The lowest node is examined first along with the links to the lowest node. It represents the frequency pattern length 1. From this, traverse path in the FP tree. This path and paths are called a conditional pattern base. Conditional pattern base is a sub-database consisting of prefix paths in the FP tree starting with the lowest node.
- ⑦ Construct a conditional FP tree which is formed by count of itemsets in the path.
- ⑧ FP are generated from conditional FP tree.





Date :

Example :

Support threshold = 50%

Confidence = 60%

Transaction	List of items
T1	I1, I2, I3
T2	I2, I3, I4
T3	I4, I5
T4	I1, I2, I4
T5	I1, I2, I3, I5
T6	I1, I2, I3, I4

Solution:

Support threshold = 50%  $\Rightarrow 0.5 * 6 = 3$

$\therefore \text{min\_sup} = 3$

1. Count each item :

item	count
I1	4
I2	5
I3	4
I4	4
I5	2

2.6

2. Sort the itemset in descending order :

item	count
I2	5
I1	4
I3	4
I4	4
I5	2

3. Build FP tree :

① considering the root node null.

② The first scan of transaction T1 : I1, I2, I3 contain

Page No. Items



Date :

- $\{I1:1\}, \{I2:1\}, \{I3:1\}$  where  $I2$  is linked as a child to root,  $I1$  is linked to  $I2$  and  $I3$  is linked to  $I1$ .
- ③  $I2: I2, I3, I4$  contain  $I2, I3$  and  $I4$  where  $I2$  is linked to root,  $I3$  is linked to  $I2$  and  $I4$  is linked to  $I3$ . But this branch would share 2 node  $I2$  node as common as it is already used in  $I1$ .
- ④ Increment the count of  $I2$  by 1 and  $I3$  is linked as a child to  $I2$ ,  $I4$  is linked as a child to  $I3$ . The count is  $\{I2:2\}, \{I3:1\}, \{I4:1\}$
- ⑤  $I3: I4, I5$ . Similarly, a new branch with  $I5$  is linked to  $I4$  as a child is created.
- ⑥  $I4: I1, I2, I4$ . The sequence will be  $I2, I1$ , and  $I4$ ,  $I2$  is already linked to root node, hence it will be incremented by 1.
- ⑦  $I5: I1, I2, I3, I5$ . The sequence will be  $I2, I1, I3$  &  $I5$ .  
Thus,  $\{I2:4\}, \{I1:3\}, \{I3:2\}, \{I5:1\}$ .
- ⑧  $I6: I1, I2, I3, I4$ . The sequence will be  $I2, I1, I3, I4$ .  
Thus,  $\{I2:5\}, \{I1:4\}, \{I3:3\}, \{I4:1\}$
4. Mining of FP-tree is summarized below:
- ① The lowest node item  $I5$  is not considered as it does not have a min support count, hence it is deleted.
- ② The next lower node is  $I4$ .  $I4$  occurs in 2 branches  $\{I2, I1, I3, I4:1\}, \{I2, I3, I4:1\}$ . Therefore considering  $I4$  as suffix the prefix paths will be  $\{I2, I1, I3:1\}, \{I2, I3:1\}$ . This forms the conditional pattern base.
- ③ The conditional pattern base is considered a transaction db, an FP tree is constructed. This will contain  $\{I2:2, I3:2\}$ .  $I1$  is not considered as it does not meet the main support count.
- ④ This path will generate all combinations of FP:  
 $\{I2: I4:2\}, \{I3, I4:2\}, \{I2, I3, I4:2\}$
- ⑤ For  $I3$ , the prefix path would be  $\{I2: I1:3\}, \{I2:1\}$ . This will generate a 2 node FP-tree  $\{I2:4, I1:3\}$  and FP generated  $\{I2, I3:4\}, \{I1, I3:3\}, \{I2, I1, I3:3\}$ .





Date :

- ⑥ For  $I_1$ , the prefix path would be :  $\{I_2:4\}$  this will generate 9 single node FP-tree :  $\{I_2:4\}$  and FP are generate  $\{I_2, I_1:4\}$

Item	conditional pattern Base	conditional FP-tree	Frequent pattern generated
$I_4$	$\{I_2, I_1, I_3:1\}$ $\{I_2, I_3:1\}$	$\{I_2:2, I_3:2\}$	$\{I_2, I_4:2\}$ , $\{I_3, I_4:2\}$ , $\{I_2, I_3, I_4:2\}$
$I_3$	$\{I_2, I_1:3\}$ , $\{I_2:1\}$	$\{I_2:4, I_1:3\}$	$\{I_2, I_3:4\}$ , $\{I_1, I_3:3\}$ $\{I_2, I_1, I_3:3\}$
$I_1$	$\{I_2:4\}$	$\{I_2:4\}$	$\{I_2, I_1:4\}$

#### Advantages:

- ① This algo needs to scan the db only twice when compared to Apriori which scans the transaction for each iteration
- ② The pairing of items is not done in this algo, and this makes it faster.
- ③ The db is stored in a compact version in memory.
- ④ It is efficient and scalable for mining both long and short frequent patterns.

#### Disadvantages:

- ① FP trees is more cumbersome and difficult to build than Apriori.
- ② It may be expensive.
- ③ When the db is large, the algo. may not fit in the shared memory.



Viva Voce : →

① Define FP-growth algorithm?

→ This algo. is an improvement to the Apriori method. A FP is generated without the need for candidate generation. FP growth algo. represents the db in the form of tree called a FP tree.

② How to construct FP tree?

→ To put it simply, an FP tree is a compressed representation of the input data. It is constructed by reading the dataset one transaction at one time and mapping each transaction onto a path in the FP-tree structure. As different transactions can have the same items, their paths may overlap.

③ Define Frequent pattern?

→ FP mining (AKA Association Rule mining) is an analytical process that finds frequent pattern, association or causal structures from data sets found in various kinds of databases such as Relational db, transactional db and other databases repositories.