

(8)

FPG Patterns, Associations and Correlations

→ FP are patterns (such as itemsets, subsequences or substructures) that appear in a data set frequently.

Ex: a set of items, such as milk and bread, that appear frequently together in a transac data set is a freq itemset.

→ FP mining searches for recurring relationships in a given data set.

→ FP mining for the discovery of interesting associations and correlations b/w itemsets in transactional and relational databases.

Motivation: user want to know

- (a) which product are often purchased together
- (b) what are subsequent purchases.

Applications: (a) Basket data analysis, catalog design, web log analysis, sale campaign analysis and DNA sequence analysis.

Importance: opens the inner properties of data sets and its foundations for many data mining tasks.

- They are
- (a) Association, correlation analysis
 - (b) Sequential, structural patterns
 - (c) pattern analysis
 - (d) classification and clustering analysis
 - (e) semantic data compression and broad app's,

$$\frac{\text{min. supp} \times \text{no. of items}}{100} = \frac{50 \times 6}{100} = 3. \text{ Threshold value}$$

no. of itemsets = 6.

confidence:

$$\text{confidence } (A \Rightarrow B) = P(B|A) = \frac{\text{support}(A \cup B)}{\text{support}(A)}$$

$$= \frac{\text{support count}(A \cup B)}{\text{support count}(A)}$$

FP mining : A Read map.

FP mining can be categorized in many diff ways

- (1) Based on the completeness of patterns to be mined (complete set of items, closed freq itemsets)
- (2) Based on the levels of dimensions of data involved in the rule.
- (3) Based on the types of values handled in the rule.
- (4) Based on the kinds of rules to be mined.
- (5) Based on the kinds of patterns to be mined.
 - can find rules at diff levels of abstraction. for ex; a set of association rules at dimension includes the following rules, buys(X, "computer") \Rightarrow buys(X, "HP-printer") → (1)
buys(X, "laptop-computer") \Rightarrow buys(X, "HP-printer") → (2)
 - here computer is a higher-level abstraction of laptop computer
 - (3) buys(X, "camp") \Rightarrow buys(X, "antivirus-software") → (3)
 - all (3) rules are single dimensional association rules b/w age(X, "30..39") \wedge income(X, "42K..48k") \Rightarrow buys(X, "high resolution TV")
 - (4) If a rule involves associations b/w the presence or absence of Boolean association rule.

Generate various kinds of rules and their interesting relationships. Also,

Market Basket Analysis:

- Freq itemset mining leads to the discovery of Associations and correlations among items in large transactional or relational data sets.
- The discovery of interesting correlation relationships among huge amounts of business transac records can help in many business decision-making processes, such as catalog design, cross-marketing and customer shopping behavior analysis.

$$\text{Support } (A \Rightarrow B) = P(A \cup B)$$

$$\text{confidence } (A \Rightarrow B) = P(B|A)$$

$$= \frac{\text{support } (A \cup B)}{\text{support } (A)} = \frac{\text{support-count}(A \cup B)}{\text{support-count}(A)} .$$

- Rules that satisfy both a minimum support threshold (min-sup) and a minimum confidence threshold (min-conf) are called strong.

Generating association rules

$I = \{I_1, I_2, I_5\}$, what are the assoc rules that can be generated from I ? The nonempty subsets of I are $\{I_1, I_2\}, \{I_1, I_5\}, \{I_2, I_5\}, \{I_1\}, \{I_2\}$ and $\{I_5\}$. The resulting assoc rules are as shown below, each listed with its confidence.

Alg: Apriori : find freq itemsets using an iterative level-wise approach based on candidate generation.

Input : $\rightarrow D$, a database of transactions;
 $\rightarrow \text{min-sup}$, the minimum support count threshold.

Output : L , freq itemsets in D .

Method :

1) $L_1 = \text{find-freq-1-itemsets}(D)$;

2) for ($k=2$; $L_{k-1} \neq \emptyset$; $k++$) {

$C_k = \text{apriori-gen}(L_{k-1})$;

 for each transc $t \in D$ // scan D for counts

$C_t = \text{subset}(C_k, t)$; // get the subsets of t that are candidates

 for each candidate $c \in C_t$

$c.\text{count}++$;

$L_k = \{ c \in C_k \mid c.\text{count} \geq \text{min-sup} \}$

} return $L = \cup_k L_k$;

Procedure apriori-gen (L_{k-1} : freq($k-1$) - itemsets)

 for each itemset $l_1 \in L_{k-1}$

 for each itemset $l_2 \in L_{k-1}$

 if $(l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2]) \wedge \dots \wedge (l_1[k-2] = l_2[k-2]) \wedge$

$(l_1[k-1] < l_2[k-1])$ then

 {

$c = l_1 \bowtie l_2$; // join step : generate candidates

 if has-infrequent-subset(c, L_{k-1}) then

 delete c ; // prune step : Remove unfruitful candidate

 else add c to C_k ;

} return C_k ;

Procedure has-infrequent-subset (c : candidate k -itemset;

L_{k-1} : freq($k-1$) - itemsets); // use prior knowledge

 for each $(k-1)$ -subset s of c

 if $s \in L_{k-1}$ then

 return TRUE;

 return FALSE;

TID	List of Items
T100	I ₁ , I ₂ , I ₅
T200	I ₂ , I ₄
T300	I ₂ , I ₃
T400	I ₁ , I ₂ , I ₄
T500	I ₁ , I ₃
T600	I ₂ , I ₃
T700	I ₁ , I ₅
T800	I ₁ , I ₂ , I ₃ , I ₅
T900	I ₁ , I ₂ , I ₃

$$\begin{aligned}
 & I_1 \wedge I_2 \Rightarrow I_5 \quad \text{conf} = 2/4 = 50\% - ① \\
 & I_1 \wedge I_5 \Rightarrow I_2 \quad \text{conf} = 2/2 = 100\% - ② \\
 & I_2 \wedge I_5 \Rightarrow I_1 \quad \text{conf} = 2/2 = 100\% - ③ \\
 & I_1 \Rightarrow I_2 \wedge I_5 \quad \text{conf} = 2/6 = 33\% - ④ \\
 & I_2 \Rightarrow I_1 \wedge I_5 \quad \text{conf} = 2/7 = 29\% - ⑤ \\
 & I_5 \Rightarrow I_1 \wedge I_2 \quad \text{conf} = 2/2 = 100\% - ⑥
 \end{aligned}$$

3.0
S can D for count of each candidate

Itemset	Sup-count
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

compare candidate support count with minimum support count = 2

itemset	sup-count
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

Generate C₂ Candidate from L₁

itemset	sup-count
I ₁ , I ₂	4
I ₁ , I ₃	4
I ₁ , I ₄	1
I ₁ , I ₅	2
I ₂ , I ₃	4
I ₂ , I ₄	2
I ₂ , I ₅	2
I ₃ , I ₄	0
I ₃ , I ₅	1
I ₄ , I ₅	0

itemset	sup-count
I ₁ , I ₂	4
I ₁ , I ₃	4
I ₁ , I ₅	2
I ₂ , I ₃	4
I ₂ , I ₄	2
I ₂ , I ₅	2
I ₃ , I ₅	2

suppose threshold value = 70%.
then rules ② ③ ⑥ satisfies association rules.

itemset	sup-count
I ₁ , I ₂ , I ₃	2
I ₁ , I ₂ , I ₅	2
I ₁ , I ₂ , I ₄	1
I ₁ , I ₃ , I ₅	1
I ₂ , I ₃ , I ₄	0
I ₂ , I ₃ , I ₅	1
I ₂ , I ₄ , I ₅	0

itemset	sup-count
I ₁ , I ₂ , I ₃	2
I ₁ , I ₂ , I ₅	2

$$\begin{aligned}
 & \text{min-sup} = 2/9 = 22\% \\
 & (\text{o}) \\
 & = \frac{\text{min-sup} \times \text{no.of itemsets}}{100} \\
 & = \frac{22\% \times 9}{100} = 1.98
 \end{aligned}$$

$$\begin{aligned}
 & \text{Confidence } (A \Rightarrow B) = P(B|A) \\
 & = \frac{\text{support-count}(A \cup B)}{\text{support-count}(A)}
 \end{aligned}$$

$\forall x \in \text{transaction}, \text{buys}(x, I_1) \wedge \text{buys}(x, I_5) \Rightarrow \text{buys}(x, I_2)$

$$[S=2, C=100\%]$$

Alg: FP-growth - Mine freq. items using an FP-tree by pattern fragment growth.

INPUT: $\rightarrow D$ a transac db

\rightarrow min-sup, the min support count threshold.

Output: The complete set of freq patterns.

Method: ① The FP-tree is constructed in the foll steps:

② Scan the transaction db D once. collect F , the set of freq items and their support counts sort F in support count desc order as L , the list of freq items.

③ Create the root of an FP-tree and label it as "null". For each

transac Trans in D do the foll

Select and sort the freq items in Trans acc to the order of L . Let the sorted freq item list in Trans be $[p|P]$, where p is the first element and P is the remaining list.

call insert-tree $([p|P], T)$, which is performed as follows.

If T has a child N such that $N.item_name = p.item_name$ then incr N 's count by 1; else create a new node N , and let its count be 1, its parent linked to be linked to T , and its node-link to the nodes with the same item-name via the node-link structure. If P is nonempty, call insert-tree (P, N) recursively.

The FP-tree is mined by calling

FP-growth (FP-tree, null), which it implemented as follows.

Procedure FP-growth (Tree, α)

- ① if Tree contains a single path P then
- ② for each combination (denoted as β) of the nodes in the path P .
- ③ generate pattern $\beta \cup \alpha$ with support_count = minimum support count of nodes in β .
- ④ else for each a_i in the header of Tree $\{ \}$
- ⑤ generate pattern $\beta = a_i \cup \alpha$ with support_count = $a_i \cdot \text{support_count}$
- ⑥ construct β 's conditional pattern base and then β 's conditional FP-tree Tree_β ,
- ⑦ if $\text{Tree}_\beta \neq \emptyset$ then
- ⑧ call FP-growth (Tree_β, β),
3

Fig: The FP-growth alg for discovering frequent itemsets without candidate generation.

Ex:

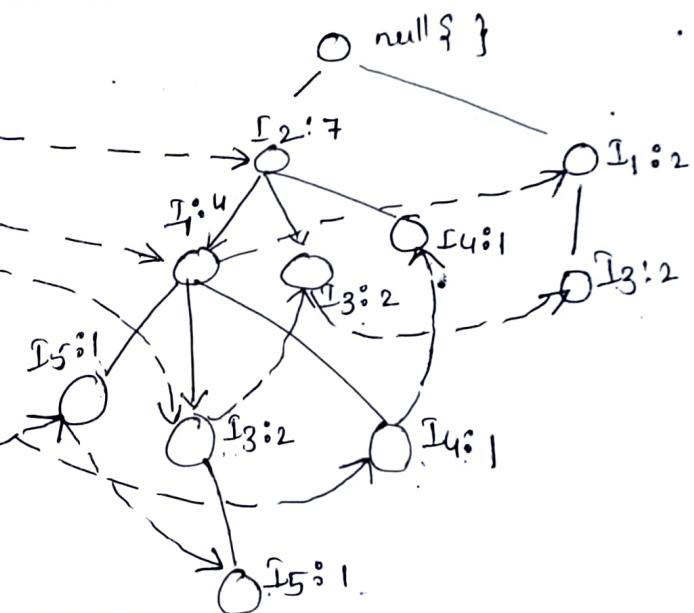
TID	list of Item-IDs
T ₁₀₀	I ₁ , I ₂ , I ₅
T ₂₀₀	I ₂ , I ₄
T ₈₀₀	I ₂ , I ₃
T ₄₀₀	I ₁ , I ₂ , I ₄
T ₅₀₀	I ₁ , I ₃
T ₆₀₀	I ₂ , I ₃
T ₇₀₀	I ₁ , I ₃
T ₈₀₀	I ₁ , I ₂ , I ₃ , I ₅
T ₉₀₀	I ₁ , I ₂ , I ₃

After descending order \Rightarrow

T ₁₀₀	I ₂ , I ₁ , I ₅
T ₂₀₀	I ₂ , I ₄
T ₃₀₀	I ₂ , I ₃
T ₄₀₀	I ₂ , I ₁ , I ₄
T ₅₀₀	I ₁ , I ₃
T ₆₀₀	I ₂ , I ₃
T ₇₀₀	I ₁ , I ₃
T ₈₀₀	I ₂ , I ₁ , I ₃ , I ₅
T ₉₀₀	I ₂ , I ₁ , I ₃

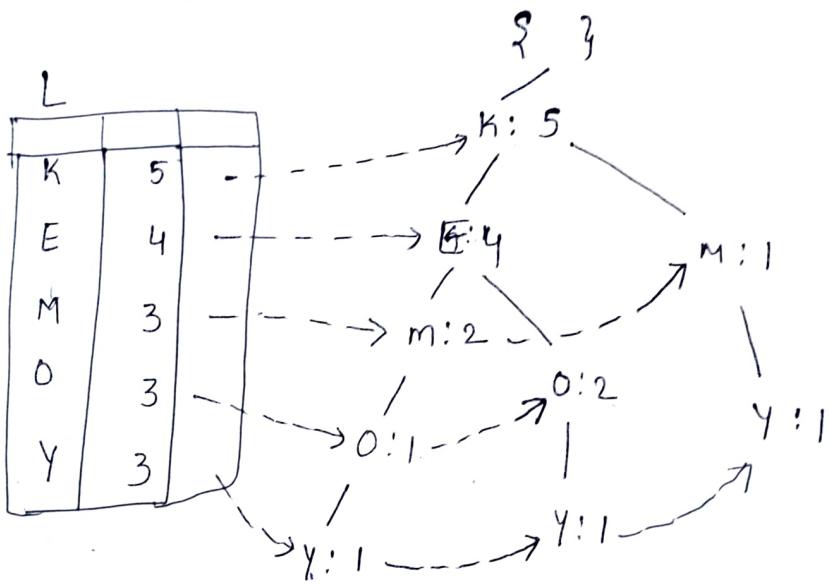
Item ID Support count Node link

Item ID	Support count	Node link
I ₂	7	-
I ₁	6	-
I ₃	6	-
I ₄	2	-
I ₅	2	-



Item	Conditional patternBase	conditionalFP-tree	Freq patterns Generated.		
I ₅	{I ₂ , I ₁ :1} {I ₂ , I ₁ , I ₃ :1}	<I ₂ :2, I ₁ :2>	{I ₂ , I ₅ :2} {I ₁ , I ₅ :2}		
I ₄	{I ₂ , I ₁ :1} {I ₂ :1}	<I ₂ :2>	{I ₂ , I ₄ :2}		
I ₃	{I ₂ , I ₁ :2}, {I ₂ :2}{I ₁ :2}	<I ₂ :4, I ₁ :2><I ₁ :2>	{I ₂ , I ₃ :4}	{I ₁ , I ₃ :4}	{I ₂ , I ₁ , I ₃ :2}
I ₁	{I ₂ :4}	<I ₂ :4>	{I ₂ , I ₁ :4}		

TID	Iluomsels	ordered itemset
T ₁₀₀	{S,M,O,N,K,E,Y}	{K, E, M, O, Y}
T ₂₀₀	{D,O,N,K,E,Y}	{K, E, O, Y}
T ₃₀₀	{M,A,K,E}	{K, E, M}
T ₄₀₀	{M,U,C,K,Y}	{K, M, Y}
T ₅₀₀	{C,O,O,K,I,E}	{K, E, O}



Items	conditional/pattern base	conditional FP tree	freq-pattern generated
Y	{K, E, M, O: 1} {K, E, O: 1} {K: M: 1}	K: 3	{K, Y: 3}
O	{K, E, M: 1} {K, E: 2}	K, E: 3	{K, O: 3} {E, O: 3}
M	{K, O: 2} {K: 1}	K: 3	{O, K, E: 3} {m, K: 3}
E	{K: 4}	K: 4	{E, K: 4}
K	-	-	-

- Improving the efficiency of Apriori
- Hash-based tech - used to reduce the size of the candidate K-itemsets, C_K , for $K > 1$.

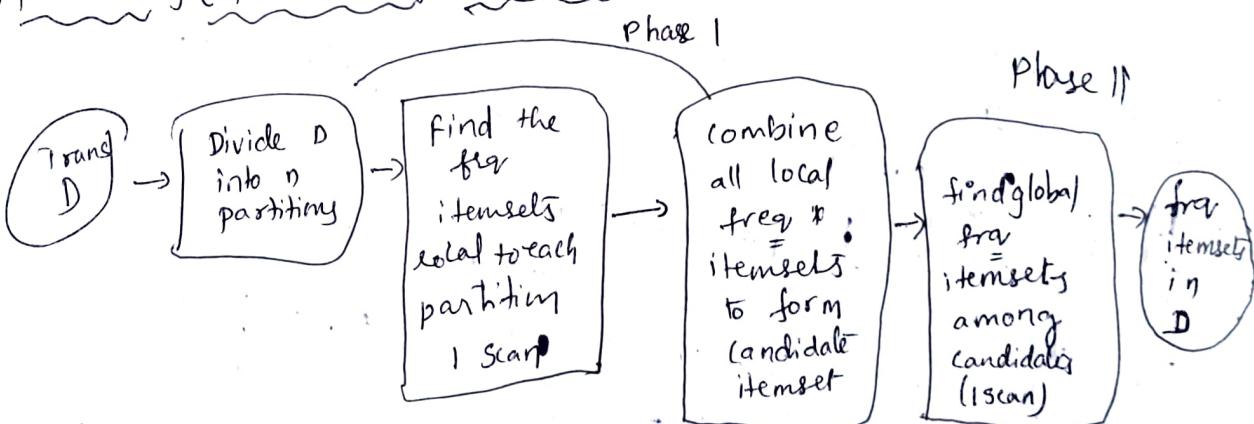
Create hashtable H_2

Using hashfunction

$$h(x,y) = ((\text{order}(x) \times 10 + (\text{order}(y))) \bmod 7)$$

bucketaddress	0	1	2	3	4	5	6
"count"	2	2	4	9	2	8	4
"contents"	$\{I_1, I_4\}$ I_3, I_5	I_1, I_5 I_1, I_5	I_2, I_3 I_2, I_3 I_2, I_3 I_2, I_3	I_2, I_4 I_2, I_4	I_2, I_5 I_2, I_5	I_1, I_2 I_1, I_2 I_1, I_2 I_1, I_2	I_1, I_3 I_1, I_3 I_1, I_3 I_1, I_3

- Transaction reduction (reducing the no. of transactions scanned in future iterations): A transac that does not contain any frequent K-itemsets cannot contain any freq ($K+1$) itemsets!
- partitioning (partitioning the data to find candidate itemsets),



- Sampling (mining on a subset of the given data): the basic idea of the Sampling approach is to pick a random sample S of the given data D , and then search for freq itemsets in S instead of D . The sample size of S is such that the search for freq itemsets in S can be done in mainmemory, and so only one scan of the transac in S is reqd overall.

Because we are searching for freq itemsets in S rather than in D , it is possible that we will miss some of the global freq itemsets.

dynamic itemset counting (adding candidate itemsets at different points during a scan);

A dynamic itemset counting tech was proposed in which the database is partitioned into blocks marked by start points. In this variation, new candidate itemsets can be added at any start point, unlike in Apriori,

which determine

1-itemset

TID	Itemset
T ₁₀₀	I ₁ I ₂ I ₅
T ₂₀₀	I ₂ I ₄
T ₃₀₀	I ₂ , I ₃
T ₄₀₀	I ₁ , I ₂ , I ₄
T ₅₀₀	I ₁ , I ₃
T ₆₀₀	I ₂ , I ₃
T ₇₀₀	I ₁ , I ₃
T ₈₀₀	I ₁ , I ₂ , I ₃ , I ₅
T ₉₀₀	I ₁ , I ₂ , I ₃

items	support count
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

Hash fun

$$h(x, y) = (\text{order of } x) * 10 + (\text{order of } y) \text{ mode } 7$$

$$I_1 = 1, I_2 = 2, I_3 = 3$$

$$I_4 = 4, I_5 = 5$$

2-itemset

Itemset	supp. count	
I ₁ , I ₂	4	(1 * 10 + 2) mode 7 = 5
I ₁ , I ₃	4	(1 * 10 + 3) mode 7 = 6
I ₁ , I ₄	1	(1 * 10 + 4) mode 7 = 0
I ₁ , I ₅	2	(1 * 10 + 5) mode 7 = 1
I ₂ , I ₃	4	(2 * 10 + 3) mode 7 = 2
I ₂ , I ₄	2	(2 * 10 + 4) mode 7 = 3
I ₂ , I ₅	2	(2 * 10 + 5) mode 7 = 4
I ₃ , I ₄	0	null
I ₃ , I ₅	1	(3 * 10 + 5) mode 7 = 0
I ₄ , I ₅	0	null

size = 7, reminder

7/12 (1)

7/13 (1)

Now consider
7/14
as this is

Bucket-
addresses

for example.

Bucket '0' contains
ignores {I₁, I₄, }
{I₃, I₅}

as count is (2)

From Association Analysis to Correlation Analysis

- A correlation measure can be used to augment the support-confidence framework for association rules. This leads to correlation rules of the form $A \Rightarrow B$ [support, confidence, correlation]

→ That is, a correlation rule is measured not only by its support and confidence but also by the correlation between items.

$A \& B$.

- Lift is a simple correlation measure that is given as follows.

→ The occurrence of itemset A is independent of the occurrence of itemset B if

$P(A \cup B) = P(A)P(B)$; otherwise, items A & B are dependent & correlated as events.

- This definition can easily be extended to more than two itemsets. The lift b/w the occurrence of A & B can be measured by computing

$$\text{lift}(A, B) = \frac{P(A \cup B)}{P(A)P(B)} \quad (1)$$

If the resulting value of eqn (1) is less than 1, then the occurrence of A is

negatively correlated with the occurrence of B.

→ If the resulting value is greater than 1, then A & B are positively correlated, meaning that the occurrence of one implies the occurrence of the other.

→ If the resulting value is equal to 1, then A & B are independent & there is no correlation b/w them.

correlation Analysis using lift

→ suppose we are interested in analyzing transaction at AllElectronics with respect¹⁵ to the purchase of computer games and videos. Let game refer to the transactions containing computer games and video refer to those containing video's.

→ of the 10,000 trans analyzed, the data of the customer trans show that 6,000 included computer games, while 7,500 included videos, and 4000 included both computer games and videos.

→ suppose that a data mining program for discovering association rules is run on the data, using a minimum support of, say, 30%, and a minimum confidence of 60%. The foll association rule is discovered:

$\text{buys}(x, \text{"computer games"}) \Rightarrow \text{buys}(x, \text{"videos"})$

[support = 40%, confidence = 66%]

①

rule ① is a strong association rule and would therefore be reported, since its support value of $\frac{4,000}{10,000} = 40\%$. & confidence value $= \frac{4,000}{6,000} = 66\%$. satisfy the min sup & min conf threshold resp.

→ In this, to help filter out misleading "strong" associations of the $A \Rightarrow B$ from the data in above ex, we need to study how the two itemsets, A & B are correlated.

→ let game refer to the trans of example that do not contain computer games, and video refers to those that do not contain videos.

- negatively correlated with the occurrence
→ If the resulting value is greater than
then A & B are positively correlated,
meaning that the occurrence of one im-
plies the occurrence of the other.
→ If the resulting value is equal to
then A & B are independent & there
is no correlation b/w them.

correlation analysis using lift

→ suppose we are interested in analyzing
transactions at AllElectronics with respect
to the purchase of computer games and
Let game refer to the transactions
containing computer games and video refer to
containing video's.

→ of the 10,000 trans = analyzed, the
show that 6,000 of the custom-
ers included computer games, while 7,
included video's, and 4000 included
computer games and video's.