

(8a) Frequent patterns and Association Rules:-

- Frequent patterns are patterns (such as itemsets, subsequences, or subsequences) that appear in a data set frequently.
For example, a set of items, such as milk and bread, that appear frequently together in a transaction data set is a frequent itemset.
- A subsequence, such as buying first a PC, then a digital camera, and then a memory card, if it occurs frequently in a shopping history database, is a frequent sequential pattern.
- A substructure can refer to different structural items, such as subgraphs, subtrees, or sublattices, which may be combined with itemsets or subsequences.
- If a substructure occurs frequently, it is called a frequent structured pattern.

- This chapter is dedicated to methods of frequent itemset mining.
- How can we find frequent itemset from large amount of data, where the data are either transactional or relational?
- How can we mine association rules in multidimensional and multi-dimensional space?
- Which association rules are the most interesting?
- How can we help to guide the mining procedure to discover interesting associations or correlations?
- How can we take advantage of user preferences as constraints to speed up the mining process?

- Market Basket Analysis : A motivational Example.
- frequent itemset mining leads to the discovery of associations and correlations among items in large transactional or relational data sets.
- A typical example of frequent itemset mining is market basket analysis.
- This process analyzes customer buying habits by finding associations between the different items that customers place in their "shopping baskets".

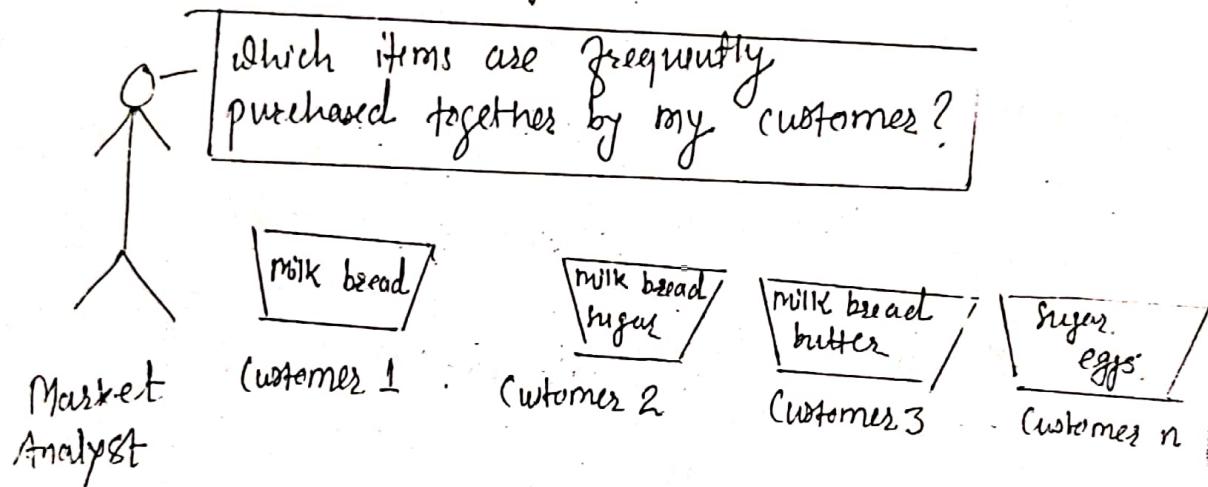


Fig. Market Basket Analysis.

- The discovery of such associations can help retailers develop marketing strategies by gaining insight into which items are frequently purchased together by customers.
- For instance, if customers are buying milk, how likely are they to also buy bread (and what kind of bread) on the same trip to the supermarket?
- Such information can lead to increased sales by helping retailers do selective marketing and plan their shelf spaces.

Market Basket Analysis:-

Suppose, if manager of an AllElectronics branch, you would like to learn more about the buying habits of your customers. Specifically, "which groups or sets of items are customers likely to purchase on a given trip to the store?

Market Basket analysis may be performed on the actual data of customer transactions at your store.

You can then use the results to plan marketing or advertising strategies or in the design of a new catalog.

For instance, market basket analysis may help you design different store layouts.

In one strategy, items that are frequently purchased together can be placed in proximity in order to further encourage the sale of such items together.

If customers who purchase computers also tend to buy antivirus software at the same time, then placing the hardware display close to the software display may help increase the sales of both items.

Market basket analysis can also help retailers plan which items to put on the sale at reduced prices.

If customers tend to purchase computers and printers together, then having a sale on printers may encourage the sale of printers as well as computers.

If we think of the universe as the set of items available at the store, then each item has a Boolean variable, signifying the presence or absence of that item.

Each basket can then be represented by a Boolean ~~variable~~ vector of values assigned to these variables.

The Boolean vectors can be analyzed for buying patterns that reflect items that are frequently associated as purchased together.

- These patterns can be represented in the form of association rules.
- for example,
 - The information that customers who purchase computers also tend to buy antivirus software at the same time it's represented in Association Rule as follows:

- $\text{Computer} \Rightarrow \text{antivirus-software}$ [support = 2%, confidence = 60%]
- Rule support and confidence are two measures of rule interestingness.
- They respectively reflect the usefulness and certainty of discovered rules.
- A support of 2% for Association rule. It means that 2% of all the transactions under analysis show that computer and antivirus software purchased together.
- A confidence of 60% means that 60% of customers who purchased a computer also bought the antivirus-software.
- Typically, association rules are considered interesting if they satisfy both a minimum support threshold and a minimum confidence threshold.
- Such thresholds can be set by users or domain experts.

- Frequent Itemsets; Closed Itemsets and Association Rules \Rightarrow
- Let $I = \{f_1, f_2, \dots, f_m\}$ be a set of items.
- Let the task relevant data, be a set of database transactions where each Transaction T is a set of items such that $T \subseteq I$.
- Let A be a set of items.
- A Transaction T is said to contain A if and only if $A \subseteq T$.
- An association rule is an implication of the form $A \Rightarrow B$.

$A \cup B \subseteq I$ and $A \cap B = \emptyset$

where I is the percentage of transactions in D that contains $A \cup B$.

This is taken to be the probability, $P(A \cup B)$.

The rule $A \Rightarrow B$ has confidence C in the transaction set D , where C is the percentage of transactions in D containing A that also contain B .

This is taken to be the conditional probability $P(B|A)$; that is,

$$\text{Support}(A \Rightarrow B) = P(A \cup B) \quad \text{--- (2)}$$

$$\text{Confidence}(A \Rightarrow B) = P(B|A). \quad \text{--- (3)}$$

$$\text{Confidence}(A \Rightarrow B) = P(B|A) = \frac{\text{Support}(A \cup B)}{\text{Support}(A)} = \frac{\text{support_count}(A \cup B)}{\text{support_count}(A)}$$

In general, association rule mining can be viewed as a two-step process:

1. Find all frequent itemset: By definition, each of these itemsets will occur at least as frequently as a predetermined minimum support count, min-sup.

2. Generate strong association rules from the frequent itemsets: By definition, these rules must satisfy minimum support and minimum confidence.

The Apriori Algorithm : Finding Frequent Itemsets Using Candidate Generation :-

- Apriori is a seminal algorithm proposed by R. Agrawal and R. Srikant in 1994 for mining frequent itemsets for Boolean association rules.
- The name of the algorithm is based on the fact that the algorithm uses prior knowledge of frequent itemset properties.
- Apriori employs an iterative approach known as a level-wise search.
- First, the set of frequent 1-itemsets is found by scanning the database to accumulate the count for each item, and collecting those items that satisfy minimum support.
- The resulting set is denoted L_1 .
- Next L_1 is used to find L_2 , the set of frequent 2-itemsets, which is used to find L_3 and so on, until no more frequent k -itemsets can be found.
- The finding of each L_k requires one full scan of the database.
- Apriori property : All nonempty subsets of a frequent itemset must also be frequent.
- How is the Apriori property used in the algorithm?
- How L_{k-1} is used to find L_k for $k \geq 2$?
- A two step process is followed, consisting of join and prune actions.
 - ⇒ The Join Step :- To find L_k , a set of candidate k -itemsets is generated by joining L_{k-1} with itself. This set of candidates denoted by C_k .
 - ⇒ The Prune Step :- C_k is a superset of L_k , that is, its members may or may not be frequent, but all of the frequent k -itemsets are included in C_k .

Transactional data for an AllElectronics branch.

ID	List of item IDs
100	I ₁ , I ₂ , I ₅
200	I ₂ , I ₄
300	I ₂ , I ₃
400	I ₁ , I ₂ , I ₄
500	I ₁ , I ₃
600	I ₂ , I ₃
700	I ₁ , I ₃
800	I ₁ , I ₂ , I ₃ , I ₅
900	I ₁ , I ₂ , I ₃

- These are nine transactions in this database. That is $|D| = 9$.
- Suppose that the minimum support count is 2.
- Here we are illustrating the Apriori Algorithm for finding frequent itemsets in D.

C ₁	L ₁
Scan D for count of each candidate.	Itemset Support Count
→	{I ₁ } 6
	{I ₂ } 7
	{I ₃ } 6
	{I ₄ } 2
	{I ₅ } 2
	Compare (candidate support count with minimum support count)
	→
	Itemset Support Count
	{I ₁ } 6
	{I ₂ } 7
	{I ₃ } 6
	{I ₄ } 2
	{I ₅ } 2

I₁ I₂
 I₁ I₃
 I₂ I₄
 I₁ I₅
 I₂ I₃
 I₂ I₄
 I₂ I₅
 I₃ I₄
 I₂ I₄

$$C_2 = L_1 \bowtie L_1$$

	Itemset	Scan D for count of each candidate	Itemset	Sup. Count
Generate C ₂ candidates from L ₁	{F ₁ , F ₂ } {F ₁ , F ₃ }		{F ₁ , F ₂ }	4 ✓
	{F ₁ , F ₄ } {F ₁ , F ₅ }		{F ₁ , F ₃ }	4 ✓
	{F ₂ , F ₃ } {F ₂ , F ₄ }		{F ₂ , F ₃ }	1 ✓
	{F ₂ , F ₅ } {F ₂ , F ₄ }		{F ₂ , F ₄ }	2 ✓
	{F ₂ , F ₃ } {F ₂ , F ₅ }		{F ₂ , F ₅ }	2 ✓
	{F ₃ , F ₄ } {F ₃ , F ₅ }		{F ₃ , F ₄ }	0 ✗
	{F ₃ , F ₅ } {F ₄ , F ₅ }		{F ₃ , F ₅ }	1 ✓
	{F ₄ , F ₅ }		{F ₄ , F ₅ }	0 ✗

$$L_2$$

Itemset	Sup. Count
{F ₁ , F ₂ }	4
{F ₁ , F ₃ }	4
{F ₁ , F ₅ }	2
{F ₂ , F ₃ }	4
{F ₂ , F ₄ }	2
{F ₂ , F ₅ }	2

$$C_3 = L_2 \bowtie L_2$$

Itemset
{F ₁ , F ₂ , F ₃ }
{F ₁ , F ₂ , F ₅ }
{F ₁ , F ₃ , F ₅ }
{F ₁ , F ₄ , F ₅ }
{F ₂ , F ₃ , F ₄ }
{F ₂ , F ₃ , F ₅ }
{F ₂ , F ₄ , F ₅ }
{F ₃ , F ₄ , F ₅ }

F₁ F₂ F₃

F₁ F₂ F₅

F₁ F₂ F₄

F₁ F₃ F₅

F₁ F₄ F₅

F₂ F₃ F₄

F₂ F₃ F₅

$$C_3$$

Itemset Sup. Count

Scan Diff. Pmt

of each

Candidate

{F₁, F₂, F₃} 2

{F₁, F₂, F₅} 2

Compare candidate

Support count with

min. support count

$$L_3$$

Itemset

Sup. Count

2

{F₁, F₂, F₅}

2

E₃

count

F₁ F₂ F₃ F₅

1

com.

104

algorithm uses $L_3 \bowtie L_3$ to generate a candidate set of 5 itemsets I_4 .

Through the join results in $\{I_1, I_2, I_3, I_5\}$, this itemset is pruned because its subset $\{I_2, I_3, I_5\}$ is not frequent with $C_4 = 0$. And the algorithm terminates, having found all of the frequent itemsets.

Algorithm: Apriori. find frequent itemset using an iterative level wise approach based on candidate generation.

Input: D, a database of transactions;
min-sup, the minimum support count threshold.

Output: L frequent itemsets in D.

Method:

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

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

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

 4) for each transaction $t \in D$ {

 5) $I_t = \text{Subset}(k, t)$;

 6) for each candidate $c \in C_t$

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

 8) }

 9) $L_k = \{c \in C_k \mid c.\text{count} > \text{min.sup}\}$

10) }

11) return $L = \cup_k L_k$;

procedure apriori_gen(L_{k-1} ; frequent($k-1$)-itemsets)

 1) for each itemset $I_1 \in L_{k-1}$

 2) for each itemset $I_2 \in L_{k-1}$

 3) if ($I_1[1] = I_2[1] \wedge I_1[2] = I_2[2] \wedge \dots \wedge (I_1[k-2] = I_2[k-2]) \wedge (I_1[k-1] \subset I_2[k-1])$)

 4) $I_1 \bowtie I_2$;

5) if has_infrequent_subset(C, L_{k-1}) then
 6) delete C ; // prune step, eliminate unfruitful candidate.
 7) else add C to L_k
 8) }
 9) return L_k ;

procedure has_infrequent_subset(C : candidate k-itemset;
 L_{k-1} : frequent $(k-1)$ -itemsets); // use prior knowledge.

1) for each $(k-1)$ -subset S of C
 2) if $S \notin L_{k-1}$ then
 3) return TRUE;
 4) return FALSE;

→ Generating Association Rules from frequent itemsets \Rightarrow

Once the frequent itemsets from transactions or a database have been found, it is straightforward to generate a strong association rules from them.

All these strong association rules satisfy both minimum support and minimum confidence.

$$\text{Confidence}(A \Rightarrow B) = P(B|A) = \frac{\text{Support-count}(A \cup B)}{\text{Support-count}(A)}$$

→ for each frequent itemset I , generate all non empty subset of I .

→ for each frequent itemset I , generate all non empty subset of I , output the rule.

→ for every non empty subset S of I , if $\frac{\text{Support-count}(I)}{\text{Support-count}(S)} > \text{min. conf}$
 $S \Rightarrow (I-S)$

→ for $I = \{F_1, F_2, F_3\}$ from above example

→ what are the association rules that can be generated from I ?

nonempty subsets of I are,
 $\{I_1, I_2\}$, $\{I_1, I_3\}$, $\{I_2, I_3\}$, $\{I_1\}$, $\{I_2\}$, and $\{I_3\}$.
The resulting association rules are,

$$I_1 \wedge I_2 \Rightarrow \text{confidence} = 2/4 = 50\%$$

$$I_1 \wedge I_3 \Rightarrow I_2 \quad \text{confidence} = \frac{\text{occurrence}}{\text{in D}}(I_1 \wedge I_3 \cup I_2)$$

$$I_2 \wedge I_3 \Rightarrow I_1 \quad \text{confidence} = \frac{\text{occurrence}}{\text{in D}}(I_2 \wedge I_3) = \frac{2}{2} = 100\%$$

$$I_1 \Rightarrow I_2 \wedge I_3 \quad \text{confidence} = \frac{2}{2} = 100\%$$

$$I_2 \Rightarrow I_1 \wedge I_3 \quad \text{confidence} = \frac{2}{6} = 33\%$$

$$I_3 \Rightarrow I_1 \wedge I_2 \quad \text{confidence} = \frac{2}{7} = 29\%$$

$$I_3 \Rightarrow I_1 \wedge I_2 \quad \text{confidence} = \frac{2}{2} = 100$$

\Rightarrow If the minimum confidence threshold is, say 70%,

\Rightarrow Then the second, third and last rules above are output, because these are the only ones generated that are strong.

Q.3) FP-growth (finding frequent itemsets without candidate generation)

- However, Apriori algorithm suffer from two non-trivial costs
- It may need to generate a huge number of candidate sets.
- It may need to repeatedly scan the database and check a large set of candidates by pattern matching.
- Frequent-pattern growth mines the complete set of frequent itemsets without candidate generation.
- FP-growth, adopts a divide-and-conquer strategy.
- First it compresses the database representing frequent items into a frequent-pattern tree, or FP-tree, which retains the itemset association information.
- It then divides the compressed database into a set of conditional databases, each associated with one frequent item or "pattern segment", and mines each database separately.

Example:- A given Transactional Data for an ALE Electronics Branch

TID List of Items-FDS

Solution \Rightarrow Frequency of Items

Item Support Count

I₁ 6

I₂ 7

I₃ 6

I₄ 2

I₅ 2

T₁ : I₁, I₂, I₅

~~2, 3, 4~~

T₂ : I₁, I₂, I₃

T₃ : I₂, I₃, I₄

T₄ : I₁, I₃, I₅

T₅ : I₁, I₂, I₃, I₄

Set of frequent items is sorted in the order of descending support count. (minimum support is 2).
After we sort by support count

I₂ . 7

I₁ . 6

I₃ . 6

I₄ . 2

I₅ . 2

3) Arrange according to polarity

TID . List of Items

T₁₀₁ I₂, I₁, I₅

T₂₀₁ I₂, I₄

T₃₀₁ I₂, I₃

T₄₀₁ I₂, I₁, I₄

T₅₀₁ I₁, I₃

T₆₀₁ I₂, I₃

T₇₀₁ I₁, I₃

T₈₀₁ I₂, I₁, I₃, I₅

T₉₀₁ I₂, I₁, I₃

4) FP-tree is then constructed

- first create the root of the tree, labeled with "null".

- Scan database D a second time.

"null"

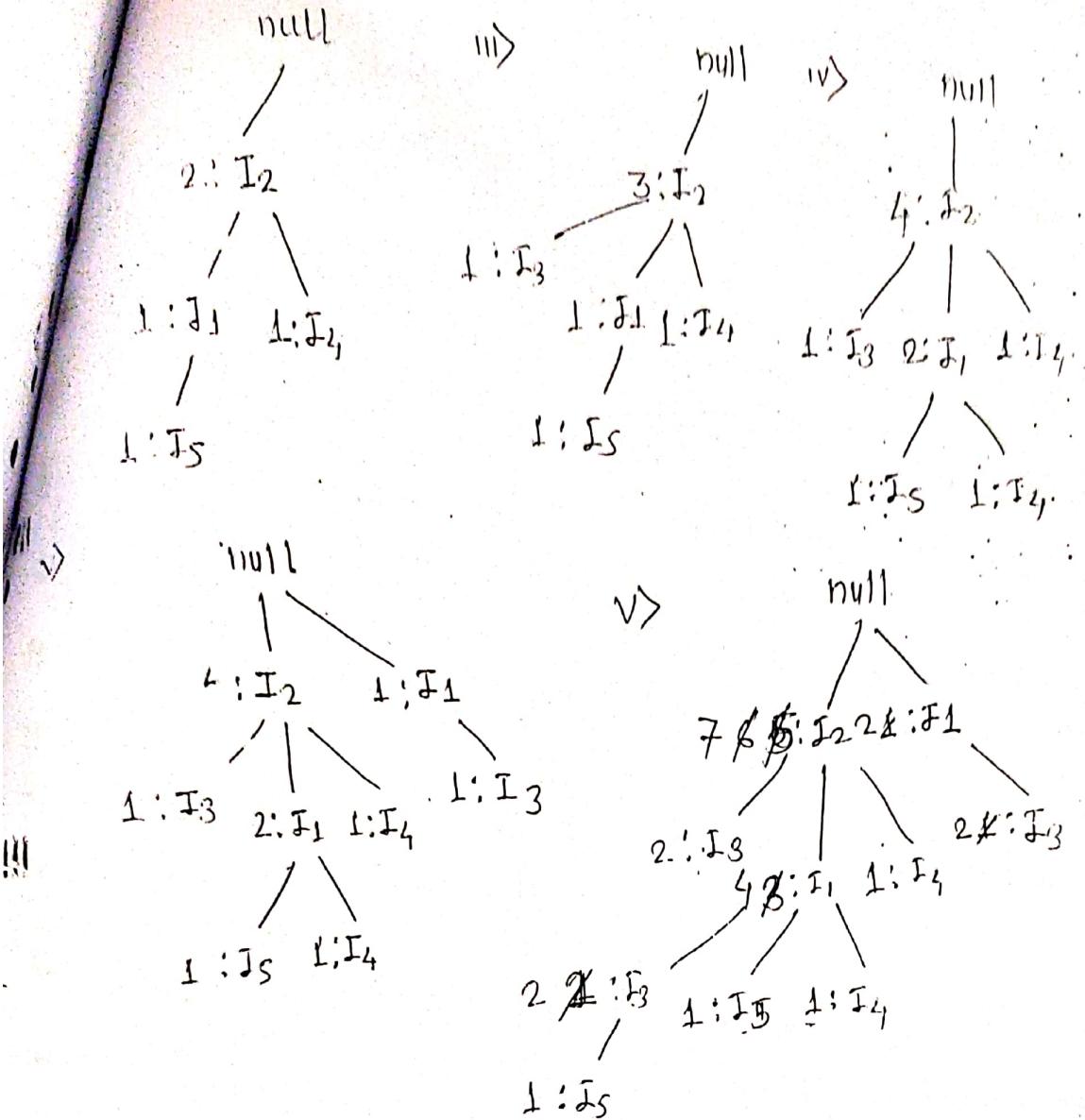
- Scan for first transaction i.e. (1)

null
|
I : I₂

|
I : I₁

|
I : I₅

109



Item	Conditional Pattern base	Conditional FPtree	Frequent Pattern Generation
I ₅	{ {I ₂ , I ₁ , I ₃ }, {I ₂ , I ₁ , I ₃ , I ₅ } }	<I ₂ :2, I ₁ :2>	{I ₂ , I ₅ :2}, {I ₁ , I ₅ :2}, {I 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 :2}, {I ₂ , I ₁ , I ₃ :2}
I ₂	{I ₂ :4}	<I ₂ :4>	{I ₂ , I ₄ :4}

Final Answer : frequent items-eds {I₁, I₂, I₃:2}
 {I₁, I₂, I₅:2}

Algorithm : FP-growth

Input : D, a transaction database;

min.sup, the minimum support count threshold.

Output : the complete set of frequent patterns.

Method :

1. The FP-tree is constructed in the following steps:

a) Scan the transaction database D once. Collect F, the set of frequent items, and their support counts. Sort F in support count descending order as L, the list of frequent items.

b) Create the root of an FP-tree, and label it as "null".

c) Sort the frequent items in L downwards according to the index of L.

d) Let the sorted frequent items list in L be [P'|P], where P' is the first element and P is the remaining list.

call insert-tree([P'|P], T),

- if T has a child N such that N.item-name = P.item-name, then increment N's count by 1. else create a new node N and let its count be 1, its parent-link 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.

2. The FP-tree is mined by calling FP-growth(FP-tree, null)

procedure FP-growth(Tree, d)

a) If Tree contains a single path P then

b) for each combination (denoted as β) of the nodes in the depth d generate pattern $\beta \cup \alpha$ with support-count = minimum m. support count of blank nodes in β .

c) else for each α_i in the headers of Tree,

generate pattern $\beta = \alpha_i \cup \alpha$ with support-count = α_i .support count

d) construct β 's conditional pattern base and then β 's conditional FP-tree $T_{\beta \cup \alpha}$;

e) if $T_{\beta \cup \alpha} \neq \emptyset$ then

f) call FP-growth($T_{\beta \cup \alpha}, \beta$); }