

Module - 3

Association Analysis

Nithin Kumar
 Aut prof
 Dept of CS&E
 VVCE, Mysore

Association Analysis :-

Association discovers the Association or Connection Among a set of items in large data sets.
 Association identifies the relationship among objects.

- Association is used for commodity Management, Advertising, Direct Marketing etc.
- For ex, A retailer can identify the products that Normally customers purchase together or even find, the customers who respond to the promotion of same kind of products.
- Consider an Example of Market basket Transaction.

TID	Items
1	{Bread, Milk}
2	{Bread, Diaper, Beer, Egg}
3	{Milk, Diaper, Beer, Cola}
4	{Bread, Milk, Diaper, Beer}
5	{Bread, Milk, Diaper, Cola}

Each row in this table corresponds to a transaction, which contains a unique identifier labeled "TID" & a set of items bought by a given customer.

- The Association Analysis is useful for discovering Purchasing relationships hidden in large data sets.
 The uncovered relationships can be represented in the form of "Association Rules" or "Sets of frequent items".

The Basic Terminology used in Association Analysis are

- * Binary Representation :- The Transaction data can be represented in a binary format

TID	Bread	Milk	Diaper	Beef	Eggs	cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

where Each row corresponds to a transaction & Each column corresponds to an item.

- An item can be treated as a binary variable whose value is 'one' if the item is present in a transaction & 'zero' otherwise.

- * Itemset :- Is a collection of one or more items.

Ex: {Milk, Bread, Diaper} = 3-Itemset

And k-Itemset is a Itemset which contains "k-items".

- * Support Count (α) :- frequency of occurrence of an Itemset.

- In the above table, {Milk, Bread, Diaper} Item appears for two times.

$$\therefore \alpha \{Milk, Bread, Diaper\} = 2$$

- * Support :- Fraction of Transaction that contains Itemset.

- The Support (s) is given by = $\frac{\text{Support Count } (\alpha)}{\text{total no of Transactions}}$

$$\therefore s \{Milk, Bread, Diaper\} = \frac{\alpha \{Milk, Bread, Diaper\}}{\text{total no of Transactions}} = 2/5 //$$

* Association Rule - An Association Rule is an Implication Expression form $X \rightarrow Y$, where X & Y are disjoint Subsets.

→ The Strength of an Association rule can be Measured in terms of

- * Support
- * Confidence.

→ Consider the foll Example

$$\{\text{Milk, diapers}\} \rightarrow \{\text{Beer}\}$$

* Support → fraction of transactions that contain both X & Y

$$S = \frac{\# \{\text{Milk, diapers, Beer}\}}{|T|} = \frac{2}{5} = 0.4 //$$

* Confidence → Measure how often Item 'y' Appear in Transaction that contain X

$$C = \frac{\# \{\text{Milk, diapers, Beer}\}}{\# \{\text{Milk, diapers}\}} = \frac{2}{3} = 0.67 //$$

→ Support is an Important Measure, Because a rule that has Very Low Support May occur Simply by chance.

Confidence on the other hand, Measures the Reliability of Interference Made by rule.

* Formulation of Association rule Mining problem - Given a Set of transactions T , the goal of Association rule Mining is to find all rules having

- * Support $>$ MinSup Threshold
- * Confidence $>$ MinConf Threshold

"MinSup" & "MinConf" are the Corresponding Support & Confidence Thresholds.

- A common strategy adopted by many association rule algorithms is to decompose the problem into two major subtasks.
 - * Frequent Itemset Generation → It is a generation of all itemsets whose support > MinSup
 - * Rule Generation → It is a generation of high confidence rules from each frequent itemset.
- * Frequent itemset Generation (Apriori Algorithm)

Apriori Algorithm

is the first association rule mining algorithm that pioneered the use of support-based pruning to systematically control the exponential growth of candidate itemsets.

- Apriori uses a "Bottom-up" approach, where frequent subsets are extended one item at a time.
- Apriori is designed to operate on database containing transactions

For Ex:- Collections of items bought by customer, or details of a website frequentation.

- * Steps to perform Apriori Algorithm
- * Step 1 :- Scan the Transaction Database to get the support of each 1-itemset, compare with "Min-Sup"
- * Step 2 :- Use Apriori property to prune the unrepresented k-itemsets from the set
- * Note :- Apriori property → Any subset of frequent itemset must be frequent

* Step 3:- Count the Support of Each Candidate by Scanning the DB

* Step 4 & 5:- Eliminate Candidates that are Infrequent, leaving only those that are frequent

Finally, Repeat the whole procedure until no new frequent itemsets are identified.

→ Consider the following Example of Frequent Itemset Generation using Apriori Algorithm

TID	Items
1	{Bread, Milk}
2	{Bread, Diaper, Beer, Egg}
3	{Milk, Diaper, Beer, Cola}
4	{Bread, Milk, Diaper, Beer}
5	{Bread, Milk, Diaper, Cola}

MinSup Count = 3

→ Firstly, Generate 1-Itemset by counting the Components

Bread	4
Diaper	4
Milk	4
Beer	3
Egg	1
Cola	2

(1-Itemset)

Eliminate the candidates that are Infrequent by Comparing with "Min Sup Count = 3" Such as "Egg & Cola"

→ Next, Generate 2-Itemset by Joining Subset of two components from 1-Itemset

$\{\text{Bread, Milk}\}$	3
$\{\text{Bread, Diaper}\}$	3
$\{\text{Bread, Beets}\}$	2
$\{\text{Milk, Diaper}\}$	3
$\{\text{Milk, Beets}\}$	2
$\{\text{Diaper, Beets}\}$	3

(2-itemset)

Get the count of Subsets by Scanning Base-DB, Eliminate the Candidates that have less count compare to "MinSup Count"

→ Generate 3-itemset by forming Subset of three Components
- to from 2-itemset.

$\{\text{Bread, Milk, Diaper}\}$	2
$\{\text{Bread, Diaper, Beets}\}$	2
$\{\text{Milk, Diaper, Beets}\}$	2
$\{\text{Bread, Milk, Beets}\}$	1

(3-itemset)

All 3-itemset are frequent, Because Count is less than Min Sup Count = 3.

* Apriori Algorithm

1. $k = 1$
2. Find all frequent 1-itemsets $F_k = \{l | l \in I \wedge \text{count}(l) \geq \text{MinSup}\}$
3. Repeat
4. $k = k + 1$
5. $F_k = \text{Apriori-gen}(F_{k-1})$ // Generate Candidate Itemsets
6. For Each Transaction $t \in T$ do
7. Generate Subset (C_k, t)
8. End for
9. $F_k = \{C | C \in C_k \wedge \text{count}(C) \geq \text{MinSup}\}$ // Extract k-frequent Itemset
10. Until $F_k = \emptyset$

II. Return F_k

- * Limitations of Apriori Algorithm → Can be very slow & the bottleneck is Candidate Generation
To compute those with Support More than MinSup, the database need to be scanned at every level. It needs $O(n^2)$. Scan. Apriori algo consumes "Huge Memory".
- * Solved problem on Frequent item Generation (Apriori) : Consider the foll Transaction

TID	Items
1	E13 E33 E43
2	E23 E33 E53
3	E13 E23 E33 E53
4	E23 E53

MinSup Count = 2

Itemset	Support
E13	2
E23	3
E33	3
E43	1
E53	3

< 2 (MinSup)

1-Itemset



Itemset	Support
E1, E3	1
E1, E3	2
E1, E5	1
E2, E3	2
E2, E5	3
E3, E5	2

< 2 (MinSup)

< 2 (MinSup)

2-Itemset



Itemset	Support
E1, 3, E3	1
E1, 3, E5	1
E1, 3, E5	1
E2, 3, E5	2

∴ Frequent Itemset = {E2, 3, E5},

* Computational Complexity of Apriori Algorithm :-
The Computational Complexity of the Apriori Algorithm can be Affected by the following factors.

* Support Threshold :- Lowering the Support Threshold ($M_{\min Sup}$) often results in More Itemsets being declared as frequent. This has an adverse Affect on Computational Complexity. of the Algorithm.

* Number of Items (Dimensionality) :- As the No of Items Increased, More Space will be needed to Store the Support Counts of Items. The Computation & I/o costs will Increase because of the Large No of Itemsets.

* Average Transaction Width :- The Average Transaction Width can be very large. This Affects the Complexity of Apriori Algorithm in two ways

- * Increase in Max Size of Frequent Itemsets
- * Increase in No of Hash Tree Traversals.

* Generation of frequent 1-Itemsets :- For Each Transaction, we need to update the Support Count for Every Item present in the Transaction.

* Candidate Generation :- To Generate Candidate k-Itemsets, pairs of frequent $(k-1)$ Itemsets are Merged. Each Merging operation requires at Most $k-2$ Equality Comparisons.

∴ Overall Cost of Merging =

$$\sum_{k=2}^w (k-2) |C_{k-1}| < \text{Cost of Merging} < \sum_{k=2}^w (k-2) |F_{k-1}|^2$$

* Support Counting :- Each Transaction of length $|t|$ produces $\binom{|t|}{k}$ items of size k . This is also the effective no of Hash tree Traversals performed for each Transaction.

* Cost for Support Counting

* Cost for Updating the Support.

* Support Counting :-

is the process of determining the frequency of occurrence for every candidate item that survives the candidate pruning of the Apriori-gen function.

→ There are two Approaches to perform Support Counting

* One Approach is to compare each transaction against every candidate item & update the support count of candidate contained in the transaction.

This Approach is computationally expensive.

* Alternative Approach is to enumerate the items contained in each transaction & use them to update the support count of their respective candidate items.

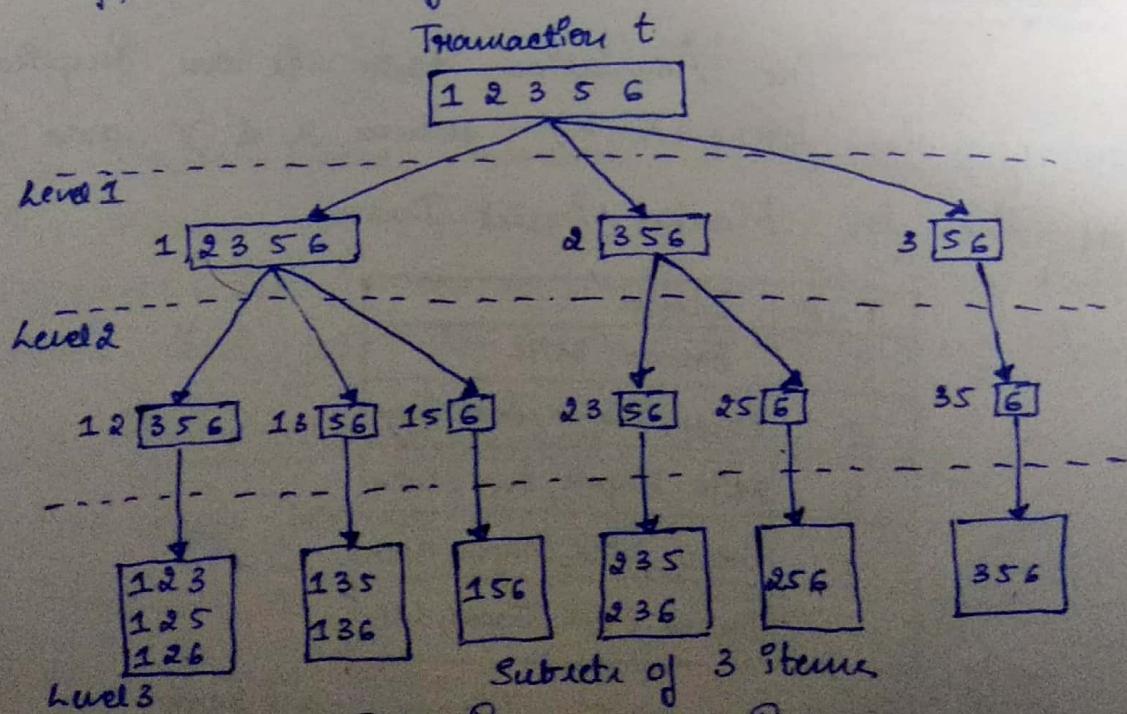


Fig:- Enumerating Subsets

- Assuming that each stemlet keeps the stems in increasing lexicographic order,
An stemlet can be enumerated by specifying the smallest stem first, followed by the larger stem.
- Consider the above fig Example, where given $t = \{1, 2, 3, 5, 6\}$. all the 3-stemlets contained in t must begin with stem 1, 2 or 3.
It is not possible to construct a 3-stemlet that begins with 5 or 6.
 - * In level 1, for instance, $1[\underline{2} \underline{3} \underline{5} \underline{6}]$ represents a 3-stemlet that begins with stem 1, followed by two more stem chosen from the set $\{2, 3, 5, 6\}$
 - * In level 2, for instance, $1\underline{2}[\underline{3} \underline{5} \underline{6}]$ corresponds to stemlets that begin with prefix (12) & are followed by stem 3, 5 or 6.
 - * Finally, the prefix structures at level 3 represent the complete set of 3-stemlets contained in t .

* Rule Generation 6-

The Association rule is an implication expression of the form $X \rightarrow Y$, where X & Y are stemlets

→ Consider the foll Market Basket Data

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Egg
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

$$\text{Ex:- Association Rule} = \{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$$

→ There are Two rule Evaluation Metrics

* Support (s) → Fraction of Transaction that contain both X & Y

* Confidence (c) → Measures how often Item Y appears in Transactions that contain X

→ For the above Mentioned Example Association rule, how to calculate Support & Confidence

$$\underbrace{\{ \text{Milk, Diaper} \}}_X \rightarrow \underbrace{\{ \text{Beer} \}}_Y$$

$$* \text{Support (s)} = \frac{\sim (X + Y)}{\text{Total no of Transactions}} = \frac{\sim \{ \text{Milk, Diaper, Beer} \}}{|T|} = \frac{2}{5} = 0.4 //$$

$$* \text{Confidence (c)} = \frac{\sim (X + Y)}{\sim (X)} = \frac{\sim \{ \text{Milk, Diaper, Beer} \}}{\sim \{ \text{Milk, Diaper} \}} = \frac{2}{3} = 0.67 //$$

→ For Given Set of Transactions T, the goal of Association rule Mining is to find all rules having

* Support \geq MinSup Threshold

* Confidence \geq MinConf Threshold

→ The Brute-Force Approach is to list all possible Association rules.

Compute the Support & Confidence for Each rule, Remove rules that fail the "MinSup" & "MinConf" Thresholds

* Rule-Generation-Algorithm (Apriori-Rule-Generation)

1. For Each frequent k-Itemset f_k , $k \geq 2$ do

2. $m = |f_{k+1}|$ // Size of rule consequent

3. If $k > m+1$ then

4. $H_{m+1} = \text{Apriori_gen}(H_m)$
5. for Each H_{m+1} do
6. if $\text{conf} \geq \text{minConf}$ then
7. Add the rule to Rule Set (H_{m+1})
8. Else
9. Delete the rule from H_{m+1}

10. End if
11. End for
12. End for

* Solved problem on Rule Generation (Apriori Rule Generation)

Consider the foll Transaction details

TID	Item
1	A, C, D
2	B, C, E
3	A, B, C, E
4	B, E

$$\text{MinSup} = 2$$

$$\text{MinConf} = 70\%$$

→ Firstly, we Should Apply Apriori Frequent Itemset Generation Algo to Generate frequent Itemset.

Itemset	S
A	2
B	3
C	3
D	1
E	3

< MinSup

1 - Itemset

Itemset	Support
A, B	1
A, C	2
A, E	1
B, C	2
B, E	3
C, E	2

2 - Itemset

< MinSup

< MinSup

3 - Itemset

Itemset	Support
A, B, C	1
A, B, E	1
A, C, E	1
B, C, E	2

< MinSup

< MinSup

< MinSup

∴ Frequent Itemset = {B, C, E}

→ Next, Apply Rule Generation procedure on obtained frequent itemset $\{B, C, E\}$ list all possible combinations.

* $B, C \rightarrow E$	Just Interchange Rule Consequent & Antecedent	$+E \rightarrow B, C$
* $C, E \rightarrow B$		$+B \rightarrow C, E$
* $B, E \rightarrow C$		$+C \rightarrow B, E$

Totally Six rules, now we should check the Confidence Value for Every individual rule.

$$* B, C \rightarrow E = \frac{\sim(B, C, E)}{\sim(B, C)} = \frac{2}{2} = 100\% > \text{MinConf} \therefore \text{Accepted} //$$

$$* C, E \rightarrow B = \frac{\sim(B, C, E)}{\sim(C, E)} = \frac{2}{2} = 100\% > \text{MinConf} \therefore \text{Accepted} //$$

$$* B, E \rightarrow C = \frac{\sim(B, C, E)}{\sim(B, E)} = \frac{2}{3} = 66.6\% < \text{MinConf} \therefore \text{Rejected}$$

$$* E \rightarrow B, C = \frac{\sim(B, C, E)}{\sim(E)} = \frac{2}{3} = 66.6\% < \text{MinConf} \therefore \text{Rejected}$$

$$* B \rightarrow C, E = \frac{\sim(B, C, E)}{\sim(B)} = \frac{2}{3} = 66.6\% < \text{MinConf} \therefore \text{Rejected}$$

$$* C \rightarrow B, E = \frac{\sim(B, C, E)}{\sim(C)} = \frac{2}{3} = 66.6\% < \text{MinConf} \therefore \text{Rejected}$$

→ Out of Six all possible combination rules, only two rules satisfied the MinConf Threshold

∴ Generated Rules are $\boxed{B, C \rightarrow E}$ & $\boxed{C, E \rightarrow B}$ //

* FP - Growth Algorithm :-

FP - Growth take a radically diff Approach to discovering Frequent Itemsets. The Algorithm does not use the "Generate-and-Test" paradigm of Apriori. Instead

It Encodes the data Set Using a Compact Data Structure called an "FP-Tree" & Extracts Frequent Itemsets directly from FP-Tree.

→ FP-Tree is an Compressed representation of Input data. Once an FP-Tree has been Constructed, it uses Recursive "Divide & Conquer" Approach to Mine frequent Itemsets.

→ FP-Growth is an Improvement of Apriori designed to Eliminate Some of the Heavy Bottlenecks in Apriori.

In FP-Tree Each node represents an Item & its Current Count & Each branch represent a different Association.

→ The FP-Growth Algorithm is divided into 5 Simple Steps as Shown below:-

* Step 1 :- The First Step is to Count all the Items in whose Transactions to get Support Count.

* Step 2 :- Next we Apply the Threshold ($MinSup$) we had set previously. & Remove Item which have Support Count less than " $MinSup$ "

* Step 3 :- Now, we Sort the list of Items with Support Count greater than " $MinSup$ " in Descending order.

* Step 4 :- Now, we Build the FP-Tree, we go through each of the Transactions & add all the Items in the order they appear in our Sorted List.

* Step 5 :- In order to get the Associations now we go through Every Branch of the Tree & only include in the Association all the nodes whose count passed the threshold ($MinSup$).

→ Consider the foll Example.

TID	Items
1	{a, b}
2	{b, c, d}
3	{a, c, d, e}
4	{a, d, e}
5	{a, b, c}
6	{a, b, c, d}
7	{a}
8	{a, b, c}
9	{a, b, d}
10	{b, c, e}

Min Support = 2

Apply FP-Growth Algorithm & Construct FP-Tree?

* Step 1 :- Count All Items In whole Transactions to get Support count

Items	Support
a	8
b	7
c	6
d	5
e	3

* Step 2 :- Apply the Threshold (MinSup) & Remove Item which power Support count less than MinSup

In the above table, all Item Support count is greater than "MinSup = 2", so all Item are considered in FP-Tree.

* Step 3 :- Arrange Items in Descending order of these Support count

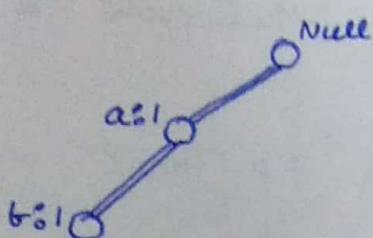
Items	Support
a	8
b	7
c	6
d	5
e	3

∴ "a, b, c, d, e" will be the FP-Tree construction pattern. (As per Descending Order)

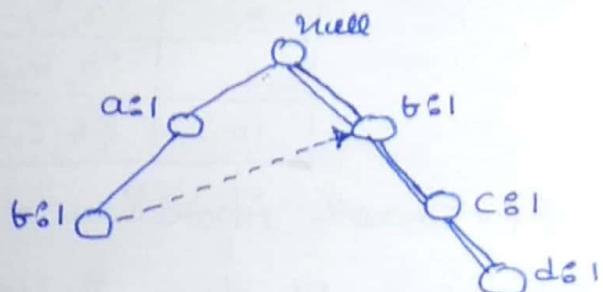
* Step 4 :- FP-Tree Construction, we go through each of Transaction with all two conditions.

- * Items whose Support Count is less than Min-Sup should be removed from Transaction
- * Descending-order Item Set (pattern) for FP-Tree Construction should be Maintained.

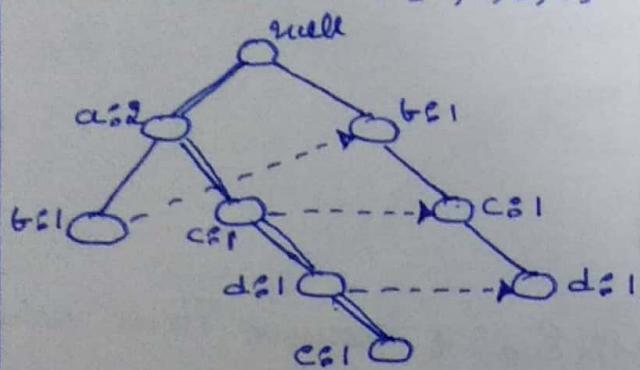
(i) Transaction - 1 $\{a, b\}$



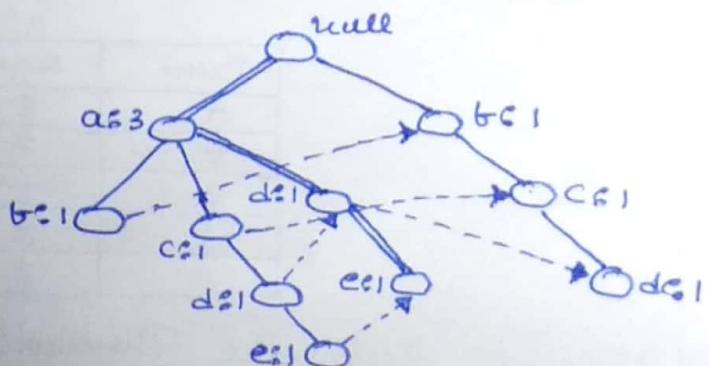
(ii) Transaction - 2 $\{b, c, d\}$



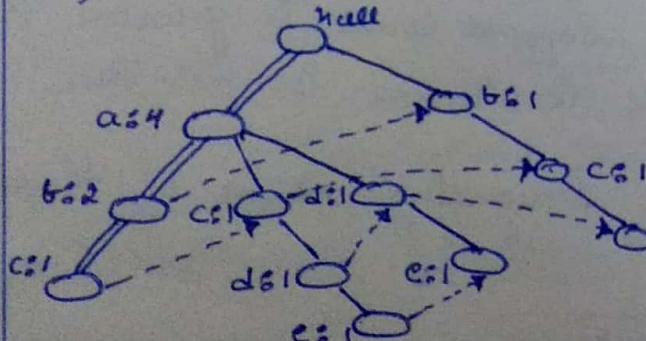
(iii) Transaction - 3 $\{a, c, d, e\}$



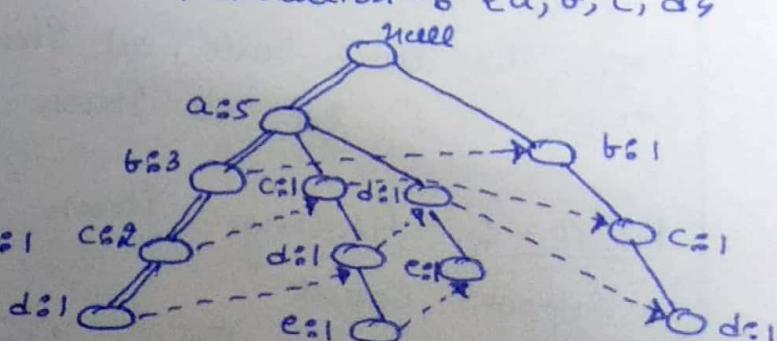
(iv) Transaction - 4 $\{a, d, e\}$



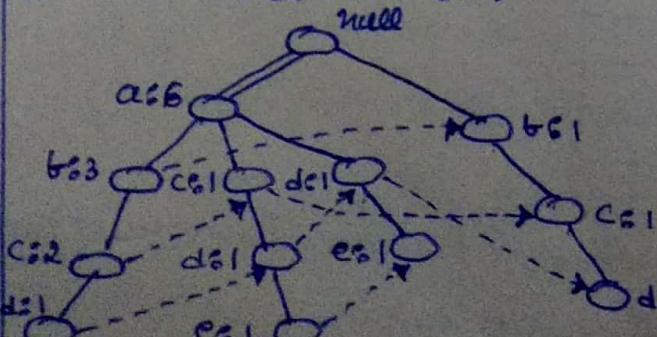
(v) Transaction - 5 $\{a, b, c\}$



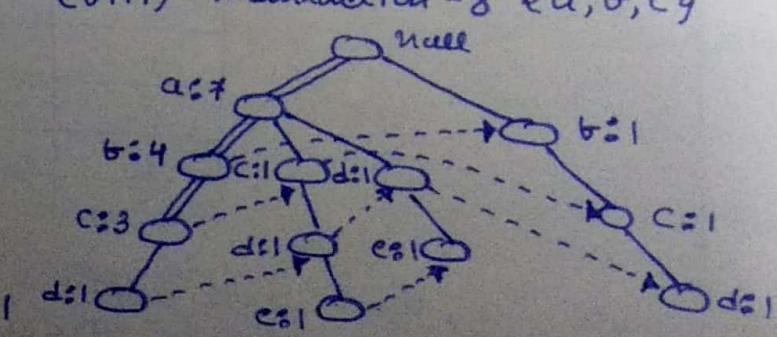
(vi) Transaction - 6 $\{a, b, c, d\}$



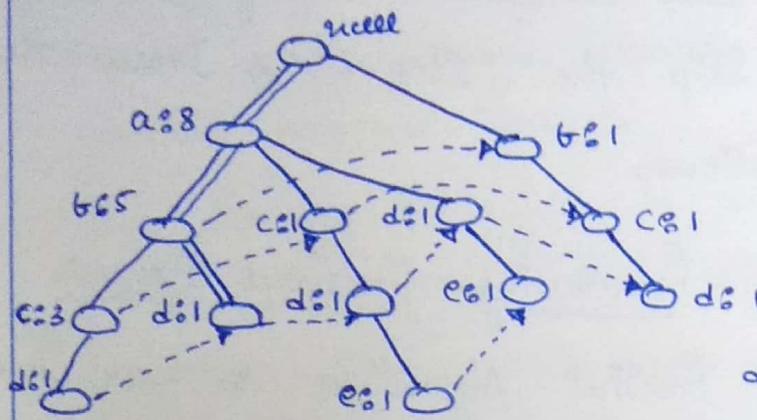
(vii) Transaction - 7 $\{a\}$



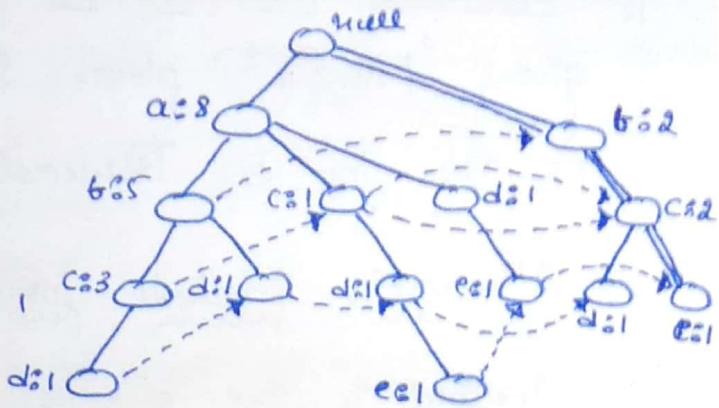
(viii) Transaction - 8 $\{a, b, c\}$



(IX) Transaction - 9 {a, b, d, e}

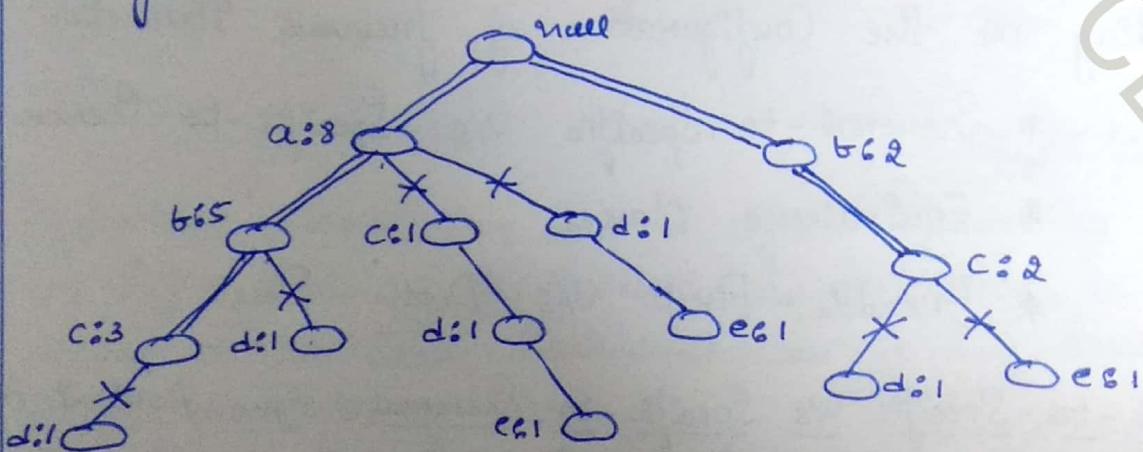


(X) Transaction - 10 {b, c, e}



(For Students Convenience, I've shown FP-Tree Construction Transaction-by-Transaction, but while solving FP-Tree problem it's better show Max Transaction Construction so that you'll get Full Marks). Note

* Step 5 - By using constructed FP-Tree compare Min Support to get Associations. $\text{MinSup} = 2$



\therefore Associations are {a, b, c} & {b, c} //

* Advantages of FP-Growth Algorithm :-

* The Biggest Advantage found in FP-Growth is the fact that Algorithm only need to read the file twice, as opposed to Apriori who reads it once for Every Iteration.

* It removes the need to calculate the pairs to be counted, which is very processing heavy.

* FP-Algo stores in Memory a compact version of DB.

* Important Note on FP-Tree Construction :- If you're not given "MinSup" please skip the "Step-2" & draw the FP-Tree for all Transactions.

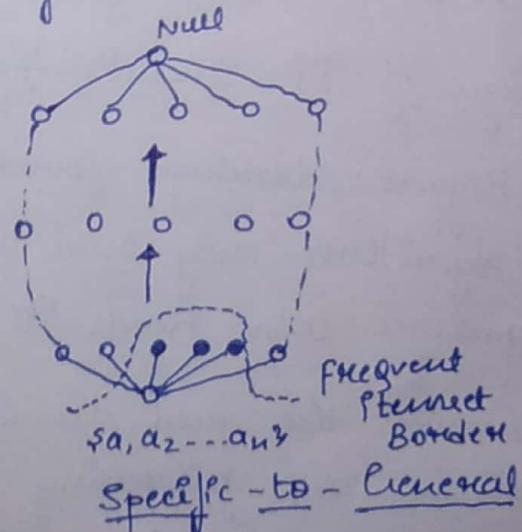
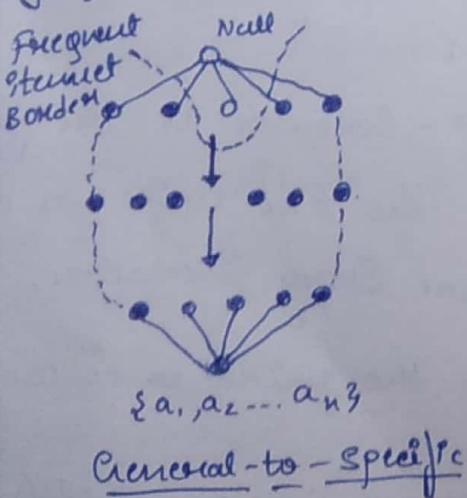
* Alternative Methods for Generating Frequent Items
Apriori is one of the earliest algorithm to have successfully addressed the combinatorial explosion of frequent itemset generation.

Despite its significant performance improvement, the algorithm still incurs considerable I/O overhead since it requires several passes over the transaction data-set.

→ Some of the search strategies are better than Apriori depending on the configuration of frequent itemsets.

- * General-to-Specific v/s Specific-to-General
- * Equivalence classes
- * Breadth-First v/s Depth-First

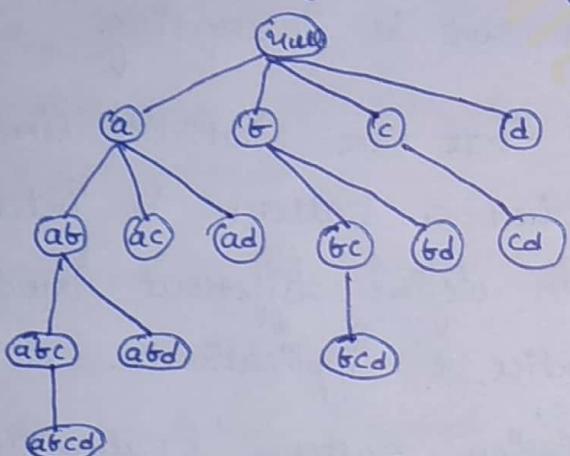
* General-to-Specific v/s Specific-to-General :- The Apriori algo uses a "general-to-specific" search strategy, where pairs of frequent (k-1) itemsets are merged to obtain k-itemsets.



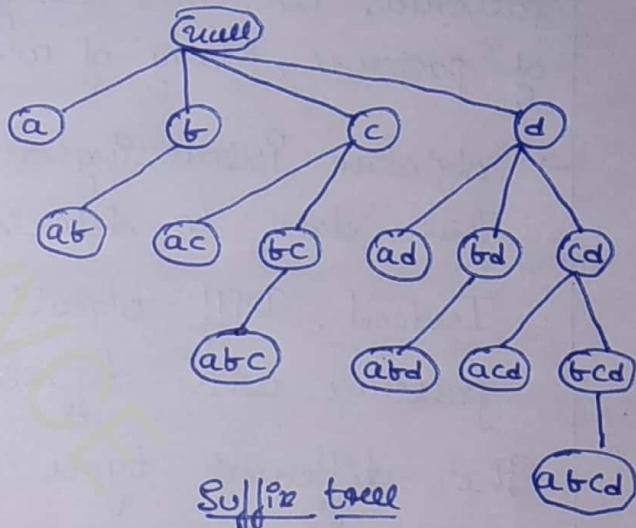
→ Alternatively, a "Specific-to-general" Search Strategy looks for more specific frequent patterns first, before finding the more general frequent patterns.

This strategy is useful to discover Maximal frequent patterns in dense transactions.

* Equivalence classes :- A frequent pattern generation algorithm searches for frequent patterns within a particular equivalence class before moving to another equivalence class.



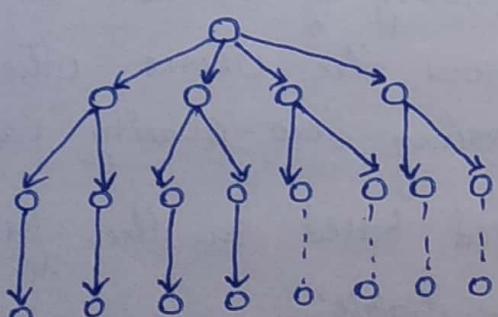
prefix Tree



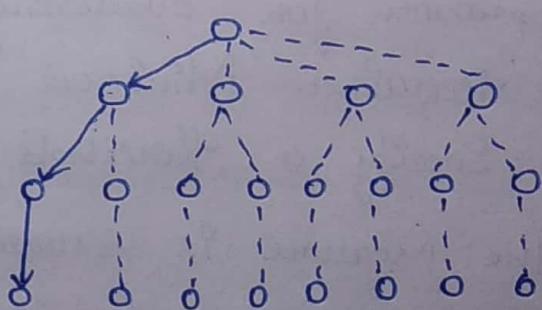
Suffix tree

Equivalence classes can also be defined according to the prefix or suffix labels of an itemset. In this case, two itemsets belong to the same equivalence class if they share a common prefix or suffix of length k.

* Breadth-First v/s Depth-First :- The Apriori algo traverses in a "Breadth-First" manner. It first discovers all the frequent 1-itemsets followed by 2-itemsets & so on.



Breadth First



Depth First

→ The "Depth-First" Approach is often used by Algorithms designed to find Maximal Frequent Patterns.

This Approach allows the frequent Pattern Border to be detected More quickly than Breadth-First Approach.

* Evaluation of Association patterns :-

The Association Analysis

Algorithms have the potential to generate a large no of patterns, we could easily end up with thousands or even millions of patterns, Many of which might not be interesting.

→ "objective Interestingness Measure" that use Statistics derived from data to determine whether a pattern is interesting. Indeed, Diff objective Measures define different Association patterns with different properties & Applications.

The different types of Association pattern Evaluation Criteria are

- * Support - Confidence Framework
- * Interest Factor
- * Correlation Analysis
- * IS Measure

* Support - Confidence Framework :- Is a data-driven Approach for Evaluating the quality of Association pattern. It requires Minimal Input from the user, other than to specify a threshold for filtering low-quality patterns.

→ The Measure is usually computed based on the frequency counts tabulated in "Contingency table"

The Contingency table is as shown below

	B	\bar{B}	
A	f_{11}	f_{10}	f_{1+}
\bar{A}	f_{01}	f_{00}	f_{0+}
	f_{+1}	f_{+0}	N

2-way Contingency table

* Limitations of Support - Confidence Framework \rightarrow Existing Association rule Mining formulation relies on the Support & Confidence Measure.

The Drawback of Support was Elimination of potentially interesting patterns & The Drawback of Confidence is More subtle.

* Interest Factor :- The High - Confidence Rules can sometimes be misleading Because the Confidence Measure Ignores the Support of Itemset Appearing in the rule - Consequent.

One way to address this problem is by Applying Metric known as "LIFT"

$$\boxed{\text{LIFT} = \frac{C(A \rightarrow B)}{S(B)}}$$

which computes the ratio between Rule Confidence & the Support of the Itemset in Rule Consequent.

* Limitations of Interest factor \rightarrow An Example from the text Mining Domain, It is reasonable to Assume that the Association between a pair of words depends on the no of Documents that contain both words.

We Expect the words "Data" & "Mining" to appear more frequently than words "Compiler" & "Mining" in Computer Articles

* Correlation Analysis :- Correlation Analysis is "Statistical Based Technique" for Analyzing Relationship between a pair of Variables.

→ For Continuous Variables, Correlation is defined using "Pearson's Correlation Co-efficient", for the Binary Variable it is given by

$$\phi = \frac{F_{11} F_{00} - F_{01} F_{10}}{\sqrt{F_{11} F_{10} F_{01} F_{00}}}$$

* Limitation of Correlation Analysis → The Drawback of Using Correlation can be seen from the word Association.

Because the ϕ -co-efficient gives equal importance to both "Co-presence" & "Co-Absence" of items in a Transaction.

* IS-Measure :- IS is an alternative Measure that has been proposed for Handling Asymmetric Binary Variables.

→ The IS Measure is defined as follows

$$* IS(A,B) = \sqrt{I(A,B) \times S(A,B)} = \frac{S(A,B)}{\sqrt{S(A) S(B)}} //$$

→ The IS Measure is large when the Pattern factor & Support of the pattern are large.

* Limitations of IS-Measure → Share a similar problem as confidence Measure that the value of the Measure can be quite large even for uncorrelated & negatively correlated patterns.