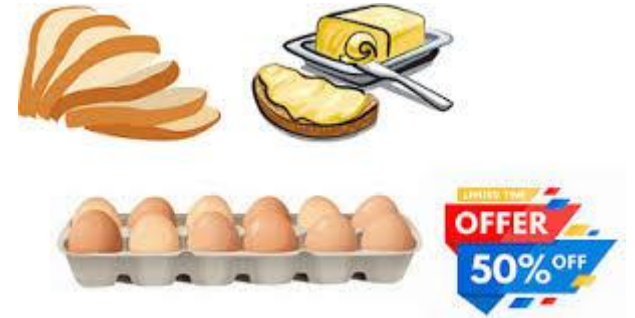


Association Rule Mining



a presentation by
Dr.B.Sivaselvan
Associate Professor, CSE
IIITD&M Kancheepuram

Association Rule Mining

- Association Rule – Implication : $X \rightarrow Y$; disjoint item-sets
- Application in Market Basket Analysis (MBA)
- Statistical Measures ; Above rule
 - Support - $P(X \cup Y)$ – Rule Significance
 - Confidence - $P(Y | X)$ – Certainty Degree
- Min Support / Confidence Strong Association Rules
- Boolean / Multidimensional / Quantitative Rules

Frequent Pattern Mining (FPM)

- Phases of Association Rule Mining
 - Frequent Item-sets Generation
 - Strong Association Rules Generation
- Apriori Algorithm –First Major contribution for FPM.
- Levelwise – Candidate generation based
- Prior Knowledge – Apriori property
- “All nonempty subsets of a frequent itemset must also be frequent.”
- Rule Generation – $s \rightarrow \{I-s\}$
 - Frequent Item-sets Generation

The Apriori Algorithm (Pseudo-Code)

C_k : Candidate itemset of size k

L_k : frequent itemset of size k

$L_1 = \{\text{frequent items}\};$

for ($k = 1; L_k \neq \emptyset; k++$) **do begin**

C_{k+1} = candidates generated from L_k ;

for each transaction t in database **do**

increment the count of all candidates in C_{k+1} that are
contained in t

L_{k+1} = candidates in C_{k+1} with min_support

end

return $\cup_k L_k$;

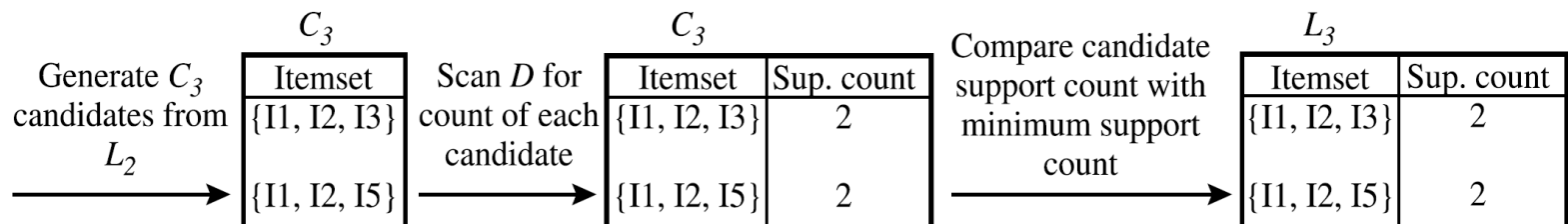
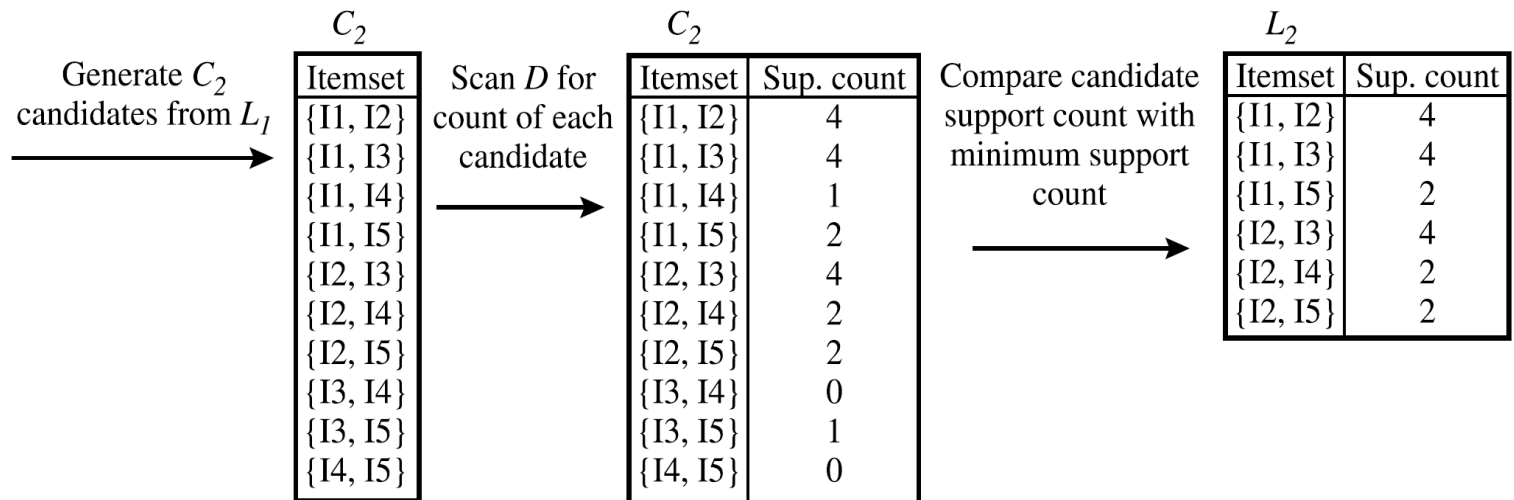
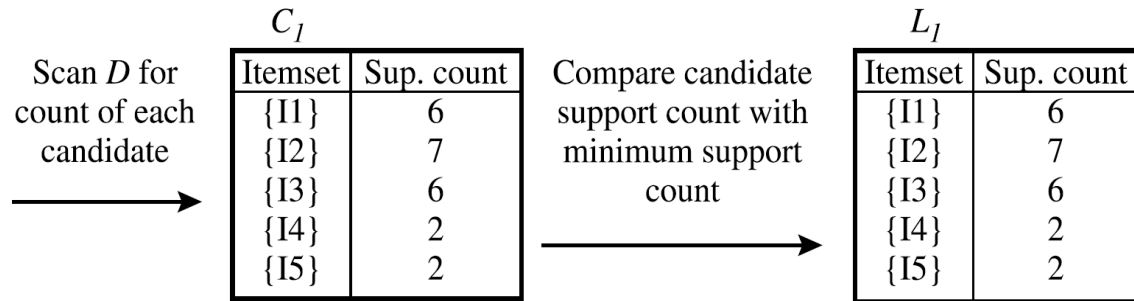
The Join and Prune Trace -

- How to generate candidates?
 - Step 1: self-joining L_k
 - Step 2: pruning
- Example of Candidate-generation
 - $L_3 = \{abc, abd, acd, ace, bcd\}$
 - Self-joining: $L_3 * L_3$
 - $abcd$ from abc and abd
 - $acde$ from acd and ace
 - Pruning:
 - $acde$ is removed because ade is not in L_3
 - $C_4 = \{abcd\}$

Apriori - Illustration

<i>TID</i>	<i>List of item_IDs</i>
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

Apriori - Illustration



- **Apriori property:** All nonempty subsets of a frequent itemset must also be frequent.
- By definition, if an itemset I does not satisfy the minimum support threshold, $\min \text{sup}$, then I is not frequent, that is, $P(I) < \min \text{sup}$. If an item A is added to the itemset I , then the resulting itemset (i.e., $I \cup A$) cannot occur more frequently than I .
Therefore, $I \cup A$ is not frequent either, that is, $P(I \cup A) < \min \text{sup}$.
- special category of properties called antimonotonicity in the sense that if a set cannot pass a test, all of its supersets will fail the same test as well.
- antimonotonicity because the property is monotonic in the context of failing a test

- $\text{confidence}(A \Rightarrow B) = P(B|A)$
- $= \text{support count}(A \cup B) / \text{support count}(A)$ For each frequent itemset I , generate all nonempty subsets of I .
- For every nonempty subset s of I , output the rule “ $s \Rightarrow (I - s)$ ” if $\text{support count}(I) / \text{support count}(s) \geq \text{min conf}$, where min conf is the minimum confidence threshold.
- $X = \{I1, I2, I5\}$
- nonempty subsets of X are $\{I1, I2\}$, $\{I1, I5\}$, $\{I2, I5\}$, $\{I1\}$, $\{I2\}$, and $\{I5\}$
- $\{I1, I2\} \Rightarrow I5$, confidence = $2/4 = 50\%$
- $\{I1, I5\} \Rightarrow I2$, confidence = $2/2 = 100\%$
- $\{I2, I5\} \Rightarrow I1$, confidence = $2/2 = 100\%$
- $I1 \Rightarrow \{I2, I5\}$, confidence = $2/6 = 33\%$
- $I2 \Rightarrow \{I1, I5\}$, confidence = $2/7 = 29\%$
- $I5 \Rightarrow \{I1, I2\}$, confidence = $2/2 = 100\%$

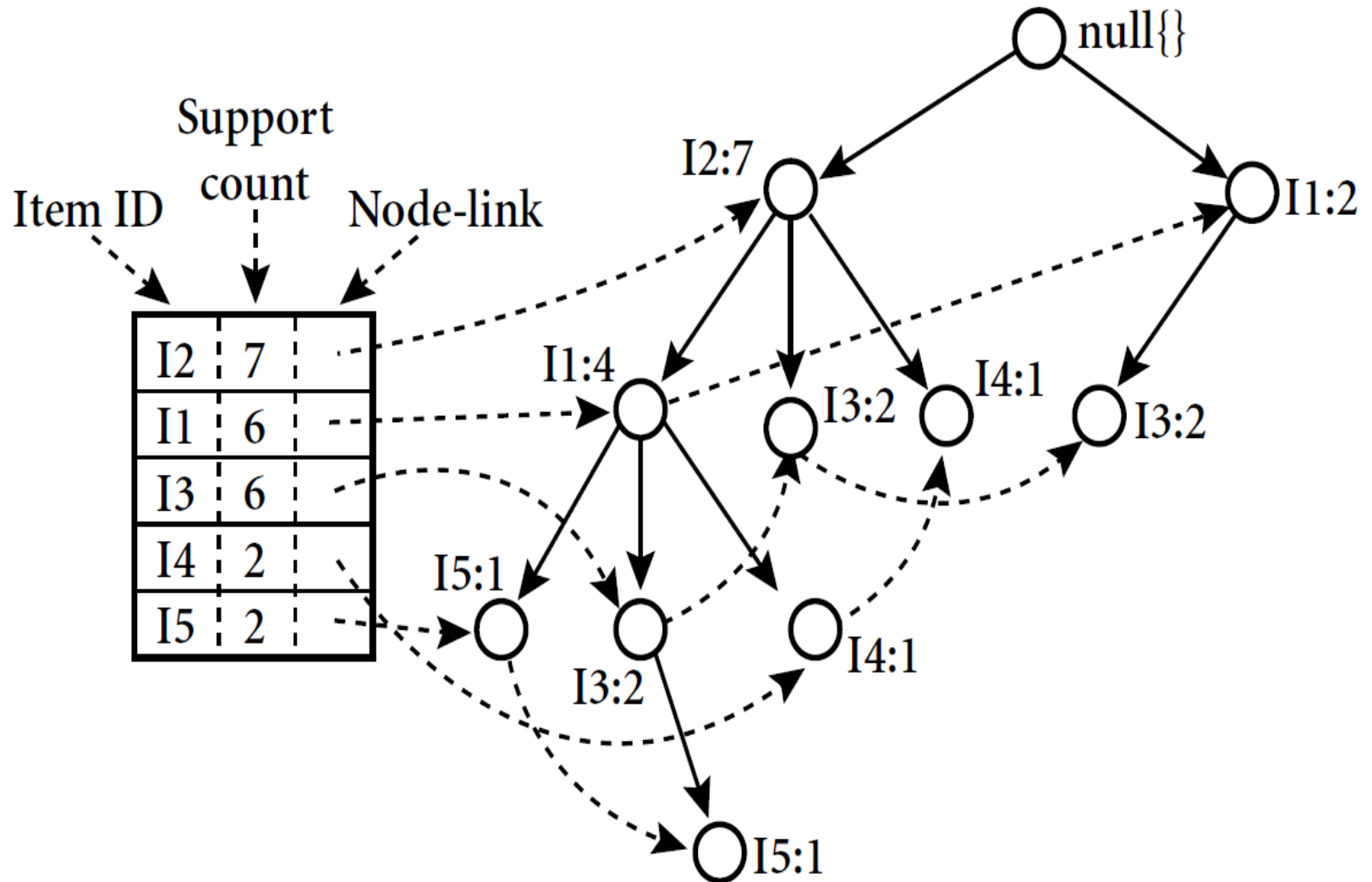
- For example, a frequent itemset of length $|I|$,
 - such as $\{a_1, a_2, \dots, a_{|I|}\}$, contains
 - $|I| C 1 = |I|$ frequent 1-itemsets: $\{a_1\}, \{a_2\}, \dots, \{a_{|I|}\}$;
 - $|I| C 2$ frequent 2-itemsets: $\{a_1, a_2\}, \{a_1, a_3\}, \dots, \{a_{|I|-1}, a_{|I|}\}$;
and so on.
 - The total number of frequent itemsets that it contains is thus
- $$|I| C 1 + |I| C 2 + \dots + |I| C |I| = 2^{|I|} - 1 \approx 1.27 \times 10^3$$

Worst-case time complexity still is exponential in $|I|$ and linear in $|D| * |I|$, but usual behavior is linear in $N = |D|$.
(detailed average-case analysis is strongly data dependent, thus difficult)

FPM – Literature

- Frequent Pattern (FP) Growth – Next Major contribution
- Improved on Apriori's Limitation – Repeated Scans of Original DB
- Overall Number of Scans – 2
- Reorders Transactions – suits MBA
- Dynamic Itemset Counting (DIC)
 - Reduced Number of Scans
 - Implication Rules – Interest and Conviction
 - Motwani et.al – Google founders

“Today, whenever you use a piece of technology, there is a good chance a little bit of Rajeev Motwani is behind it”



Item	Conditional Pattern Base	Conditional FP-tree	Frequent Patterns Generated
I5	$\{\{I2, I1: 1\}, \{I2, I1, I3: 1\}\}$	$\langle I2: 2, I1: 2 \rangle$	$\{I2, I5: 2\}, \{I1, I5: 2\}, \{I2, I1, I5: 2\}$
I4	$\{\{I2, I1: 1\}, \{I2: 1\}\}$	$\langle I2: 2 \rangle$	$\{I2, I4: 2\}$
I3	$\{\{I2, I1: 2\}, \{I2: 2\}, \{I1: 2\}\}$	$\langle I2: 4, I1: 2 \rangle, \langle I1: 2 \rangle$	$\{I2, I3: 4\}, \{I1, I3: 4\}, \{I2, I1, I3: 2\}$
I1	$\{\{I2: 4\}\}$	$\langle I2: 4 \rangle$	$\{I2, I1: 4\}$