

Data Mining and Web Algorithms

Course Code: 15B22CI621

Credits: 4 [3+ 1]

Frequent Pattern Mining: An Example

Given a transaction database DB and a minimum support threshold ξ , find all frequent patterns (item sets) with support no less than ξ .

| | | | |
|--------|-----|------------|------------------------------|
| Input: | DB: | <u>TID</u> | <u>Items bought</u> |
| | | 100 | $\{f, a, c, d, g, i, m, p\}$ |
| | | 200 | $\{a, b, c, f, l, m, o\}$ |
| | | 300 | $\{b, f, h, j, o\}$ |
| | | 400 | $\{b, c, k, s, p\}$ |
| | | 500 | $\{a, f, c, e, l, p, m, n\}$ |

Minimum support: $\xi = 3$

Output: all frequent patterns, i.e., $f, a, \dots, fa, fac, fam, fm, am \dots$

Problem Statement: How to **efficiently** find all frequent patterns?

Apriori

- Main Steps of Apriori Algorithm:

- Use frequent $(k - 1)$ -itemsets (L_{k-1}) to generate **candidates** of frequent k -itemsets C_k
- Scan database and count each pattern in C_k , get frequent k -itemsets (L_k).

Candidate
Generation



Candidate
Test



- E.g. ,

| <u>TID</u> | <u>Items bought</u> | <u>Apriori</u> | <u>iteration</u> |
|------------|--------------------------|----------------|-----------------------------------|
| 100 | {f, a, c, d, g, i, m, p} | C1 | f,a,c,d,g,i,m,p,l,o,h,j,k,s,b,e,n |
| 200 | {a, b, c, f, l, m, o} | L1 | f, a, c, m, b, p |
| 300 | {b, f, h, j, o} | C2 | fa, fc, fm, fp, ac, am, ...bp |
| 400 | {b, c, k, s, p} | L2 | fa, fc, fm, ... |
| 500 | {a, f, c, e, l, p, m, n} | ... | |

Performance Bottlenecks of Apriori

- Bottlenecks of *Apriori*: candidate generation
 - Generate huge candidate sets.
 - **Candidate Test** incur multiple scans of database.

Overview of FP-Growth: Ideas

- Compress a large database into a compact, *Frequent-Pattern tree* (FP-tree) structure
 - highly compacted, but complete for frequent pattern mining
 - avoid costly repeated database scans
- Develop an efficient, FP-tree-based frequent pattern mining method (FP-growth)
 - A divide-and-conquer methodology: decompose mining tasks into smaller ones
 - Avoid candidate generation: sub-database test only.

FP-Tree

FP-tree: Construction and Design

Construct FP-tree

Two Steps:

1. Scan the transaction DB for the first time, find frequent items (single item patterns) and order them into a list **L** in frequency descending order.
2. For each transaction, order its frequent items according to the order in **L**; Scan DB the second time, construct FP-tree by putting each **frequency ordered transaction** onto it.

FP-tree Example: step 1

Step 1: **Scan DB** for the first time to generate **L**

| <i>TID</i> | <i>Items bought</i> |
|------------|-------------------------------------------------------------------------------------------|
| 100 | { <i>f</i> , <i>a</i> , <i>c</i> , <i>d</i> , <i>g</i> , <i>i</i> , <i>m</i> , <i>p</i> } |
| 200 | { <i>a</i> , <i>b</i> , <i>c</i> , <i>f</i> , <i>l</i> , <i>m</i> , <i>o</i> } |
| 300 | { <i>b</i> , <i>f</i> , <i>h</i> , <i>j</i> , <i>o</i> } |
| 400 | { <i>b</i> , <i>c</i> , <i>k</i> , <i>s</i> , <i>p</i> } |
| 500 | { <i>a</i> , <i>f</i> , <i>c</i> , <i>e</i> , <i>l</i> , <i>p</i> , <i>m</i> , <i>n</i> } |



| <i>Item</i> | <i>frequency</i> |
|-------------|------------------|
| <i>f</i> | 4 |
| <i>c</i> | 4 |
| <i>a</i> | 3 |
| <i>b</i> | 3 |
| <i>m</i> | 3 |
| <i>p</i> | 3 |



By-Product of First Scan
of Database

FP-tree Example: step 2

Step 2: scan the DB for the second time, order frequent items in each transaction

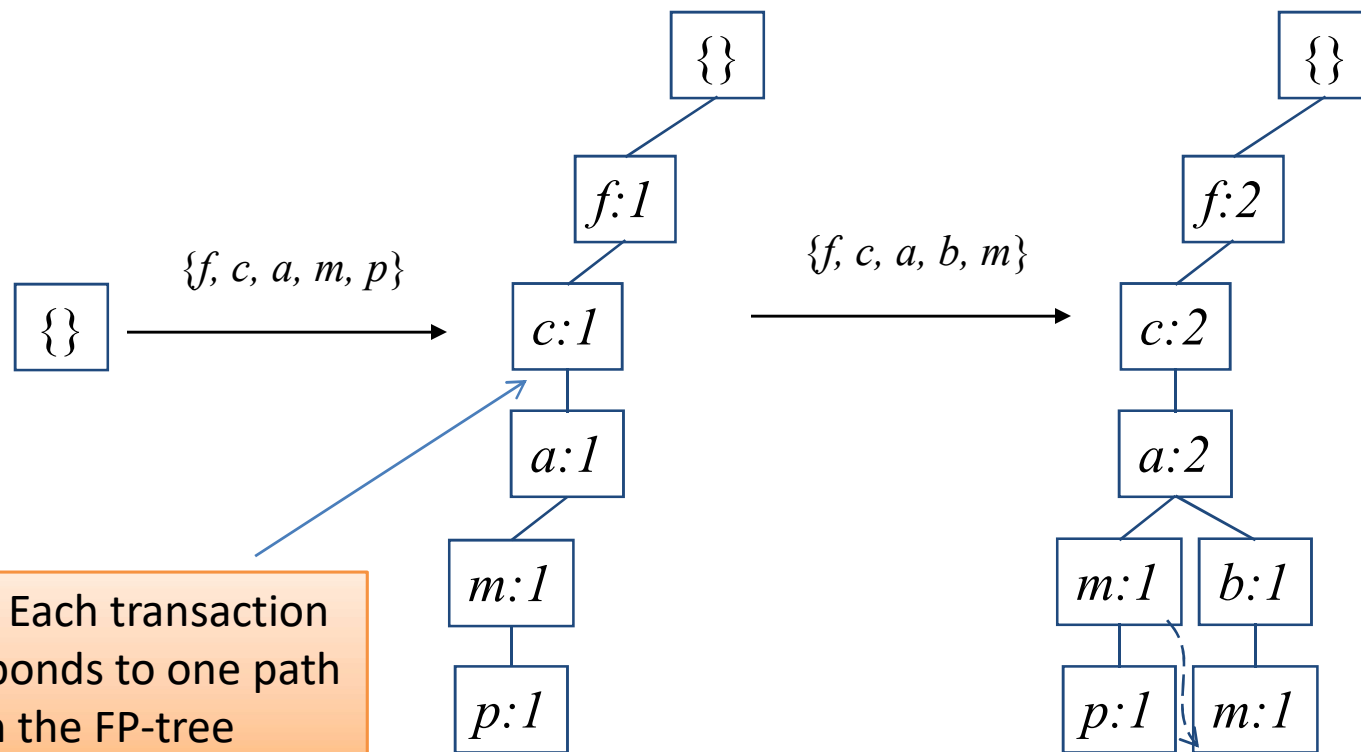
| <i>TID</i> | <i>Items bought</i> | <i>(ordered) frequent items</i> |
|------------|--------------------------|---------------------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o} | {f, b} |
| 400 | {b, c, k, s, p} | {c, b, p} |
| 500 | {a, f, c, e, l, p, m, n} | {f, c, a, m, p} |

| <i>Item</i> | <i>frequency</i> |
|-------------|------------------|
| f | 4 |
| c | 4 |
| a | 3 |
| b | 3 |
| m | 3 |
| p | 3 |

FP-tree Example: step 2

Step 2: construct FP-tree

| <i>TID</i> | <i>Items bought</i> | <i>(ordered)</i> |
|------------|--------------------------|------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o} | {f, b} |
| 400 | {b, c, k, s, p} | {c, b, p} |
| 500 | {a, f, c, e, l, p, m, n} | {f, c, a, m, p} |

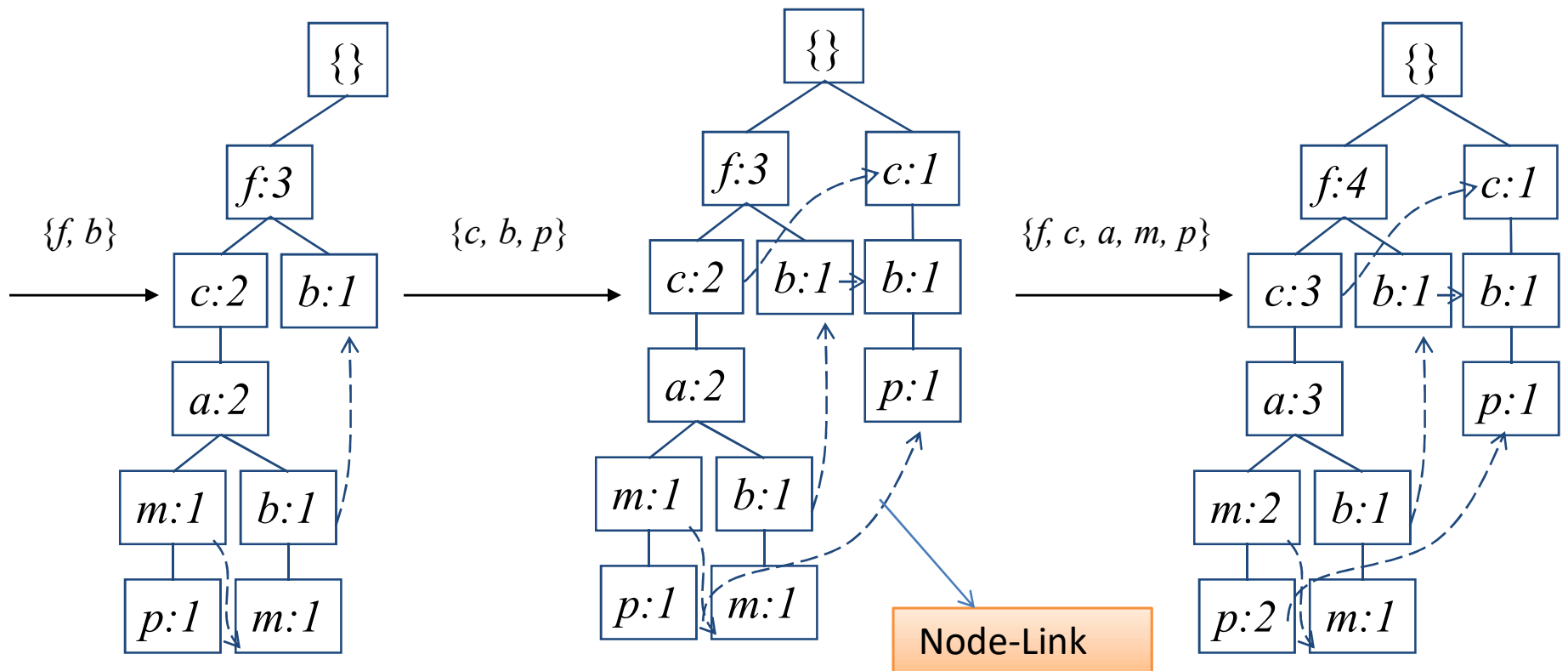


NOTE: Each transaction corresponds to one path in the FP-tree

FP-tree Example: step 2

Step 2: construct FP-tree

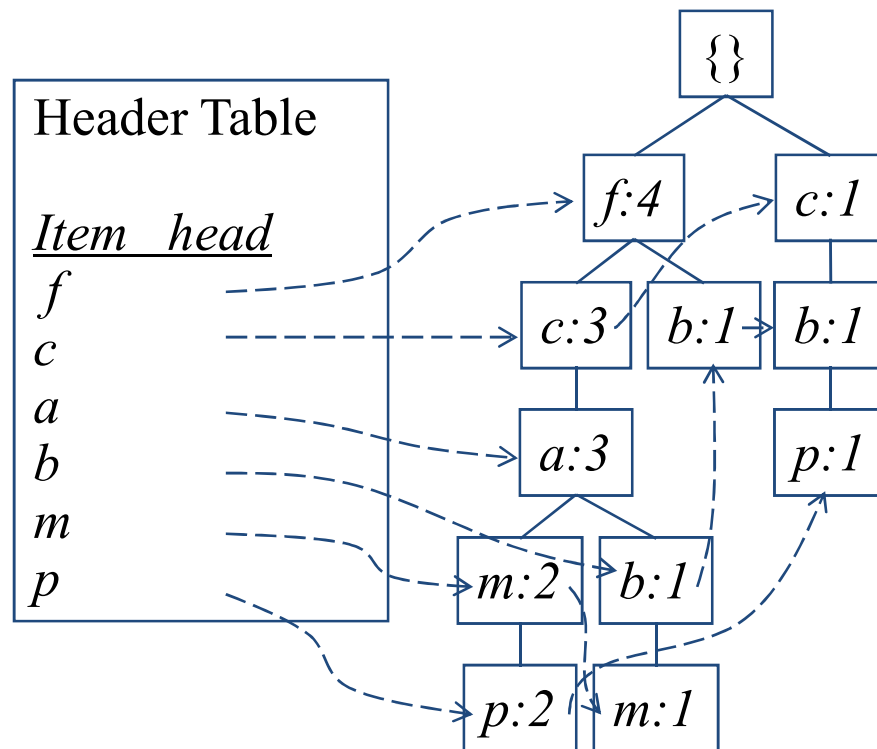
| TID | Items bought | (ordered) |
|-----|--------------------------|------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o} | {f, b} |
| 400 | {b, c, k, s, p} | {c, b, p} |
| 500 | {a, f, c, e, l, p, m, n} | {f, c, a, m, p} |



Construction Example

Final FP-tree

| <u>TID</u> | <u>Items bought</u> | <u>(ordered)</u> |
|------------|--------------------------|------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o} | {f, b} |
| 400 | {b, c, k, s, p} | {c, b, p} |
| 500 | {a, f, c, e, l, p, m, n} | {f, c, a, m, p} |



Advantages of the FP-tree Structure

- The most significant advantage of the FP-tree
 - Scan the DB only twice and twice only.
 - Completeness:
 - the FP-tree contains all the information related to mining frequent patterns (given the min-support threshold). Why?
 - Compactness:
 - The size of the tree is bounded by the occurrences of frequent items
 - The height of the tree is bounded by the maximum number of items in a transaction
-

FP-Growth

FP-growth:
Mining Frequent Patterns
Using FP-tree

Properties of FP-Tree

- Node-link property
 - For any frequent item a_i , all the possible frequent patterns that contain a_i can be obtained by following a_i 's node-links, starting from a_i 's head in the FP-tree header.
- Prefix path property
 - To calculate the frequent patterns for a node a_i in a path P , only the prefix sub-path of a_i in P need to be accumulated, and its frequency count should carry the same count as node a_i .

Principles of FP-Growth

- Pattern growth property
 - Let α be a frequent itemset in DB, B be α 's conditional pattern base, and β be an itemset in B . Then $\alpha \cup \beta$ is a frequent itemset in DB iff β is frequent in B .
 - Is “*fcabm*” a frequent pattern?
 - “*fcab*” is a branch of *m*'s conditional pattern base
 - “*b*” is **NOT** frequent in transactions containing “*fcab*”
 - “*bm*” is **NOT** a frequent itemset.
-

Mining Frequent Patterns Using FP-tree

- General idea (divide-and-conquer)

Recursively grow frequent patterns using the FP-tree: looking for shorter ones recursively and then concatenating the suffix:

- For each frequent item, construct its **conditional pattern base**, and then its **conditional FP-tree**;
- Repeat the process on each newly created conditional FP-tree until the resulting FP-tree is empty, or it contains only one **path** (single path will generate all the combinations of its sub-paths, each of which is a frequent pattern)

Major Steps

Starting the processing from the end of list **L**:

Step 1:

Construct **conditional pattern base** for each item in the header table

Step 2

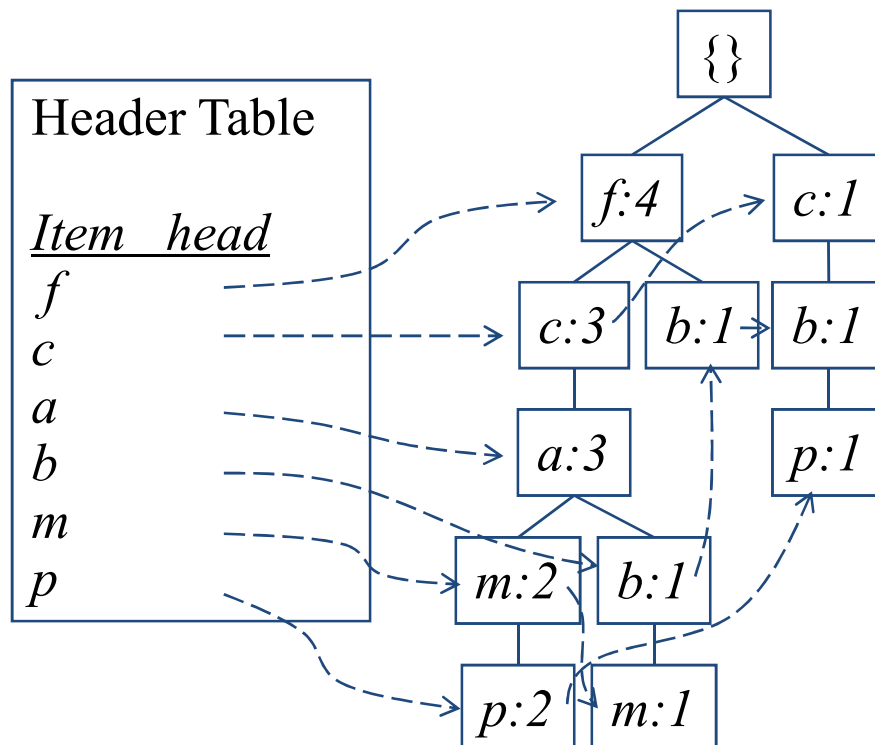
Construct **conditional FP-tree** from each conditional pattern base

Step 3

Recursively mine conditional FP-trees and grow frequent patterns obtained so far. If the conditional FP-tree contains a **single path**, simply enumerate all the patterns

Step 1: Construct Conditional Pattern Base

- Starting at the bottom of frequent-item header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item
- Accumulate all of **transformed prefix paths** of that item to form a **conditional pattern base**



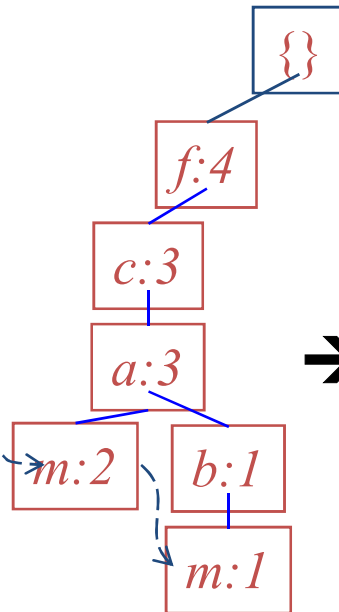
Conditional pattern bases

| <i>item</i> | <i>cond. pattern base</i> |
|-------------|---------------------------|
| <i>p</i> | <i>fcam:2, cb:1</i> |
| <i>m</i> | <i>fca:2, fcab:1</i> |
| <i>b</i> | <i>fca:1, f:1, c:1</i> |
| <i>a</i> | <i>fc:3</i> |
| <i>c</i> | <i>f:3</i> |
| <i>f</i> | <i>{}</i> |

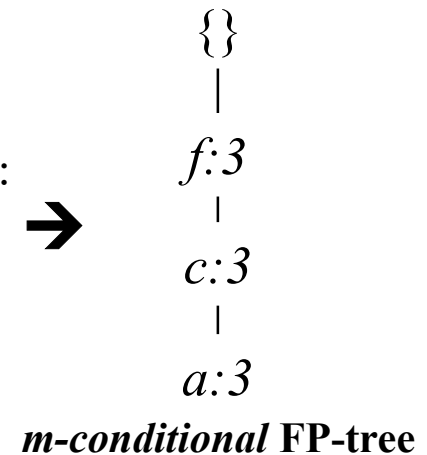
Step 2: Construct Conditional FP-tree

- For each pattern base
 - Accumulate the count for each item in the base
 - Construct the **conditional FP-tree** for the frequent items of the pattern base

| Header Table | |
|--------------|-------------|
| <i>Item</i> | <i>head</i> |
| <i>f</i> | 4 |
| <i>c</i> | 4 |
| <i>a</i> | 3 |
| <i>b</i> | 3 |
| <i>m</i> | 3 |
| <i>p</i> | 3 |



\rightarrow **m- cond. pattern base:**
 $fca:2, fcab:1$



| Conditional pattern bases | |
|---------------------------|---------------------------|
| <i>item</i> | <i>cond. pattern base</i> |
| <i>p</i> | $fcam:2, cb:1$ |
| <i>m</i> | $fca:2, fcab:1$ |
| <i>b</i> | $fca:1, f:1, c:1$ |
| <i>a</i> | $fc:3$ |
| <i>c</i> | $f:3$ |
| <i>f</i> | $\{\}$ |

Conditional Pattern Bases and Conditional FP-Tree

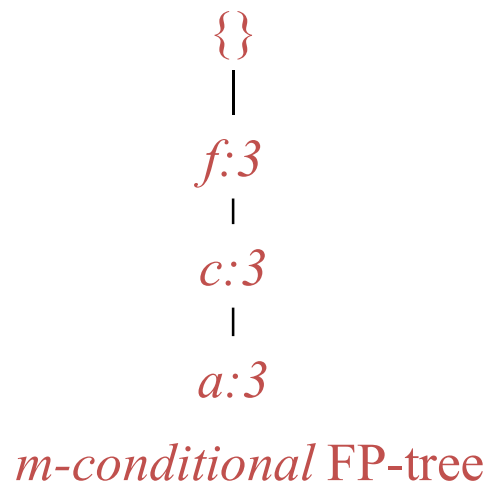
| <i>Conditional pattern bases</i> | |
|----------------------------------|---------------------------|
| <i>item</i> | <i>cond. pattern base</i> |
| <i>p</i> | <i>fcam:2, cb:1</i> |
| <i>m</i> | <i>fca:2, fcab:1</i> |
| <i>b</i> | <i>fca:1, f:1, c:1</i> |
| <i>a</i> | <i>fc:3</i> |
| <i>c</i> | <i>f:3</i> |
| <i>f</i> | <i>{}</i> |

| Item | Conditional pattern base | Conditional FP-tree |
|------|----------------------------|---------------------|
| p | {(fcam:2), (cb:1)} | {(c:3)} p |
| m | {(fca:2), (fcab:1)} | {(f:3, c:3, a:3)} m |
| b | {(fca:1), (f:1), (c:1)} | Empty |
| a | {(fc:3)} | {(f:3, c:3)} a |
| c | {(f:3)} | {(f:3)} c |
| f | Empty | Empty |

order of L

Single FP-tree Path Generation

- Suppose an FP-tree T has a single path P. The complete set of frequent pattern of T can be generated by enumeration of all the combinations of the sub-paths of P



All frequent patterns concerning *m*:
combination of {f, c, a} and *m*

m,
fm, *cm*, *am*,
fcm, *fam*, *cam*,
fcam

Example : FP Growth

| TID | List of Items |
|------|----------------|
| T100 | I1, I2, I5 |
| T101 | I2, I4 |
| T102 | I2, I3 |
| T103 | I1, I2, I4 |
| T104 | I1, I3 |
| T105 | I2, I3 |
| T106 | I1, I3 |
| T107 | I1, I2, I3, I5 |
| T108 | I1, I2, I3 |

- Consider a database, D , consisting of 9 transactions.
- Let support count =2

| Itemset | Sup.Count |
|---------|-----------|
| {I1} | 6 |
| {I2} | 7 |
| {I3} | 6 |
| {I4} | 2 |
| {I5} | 2 |

Example : FP Growth cont...

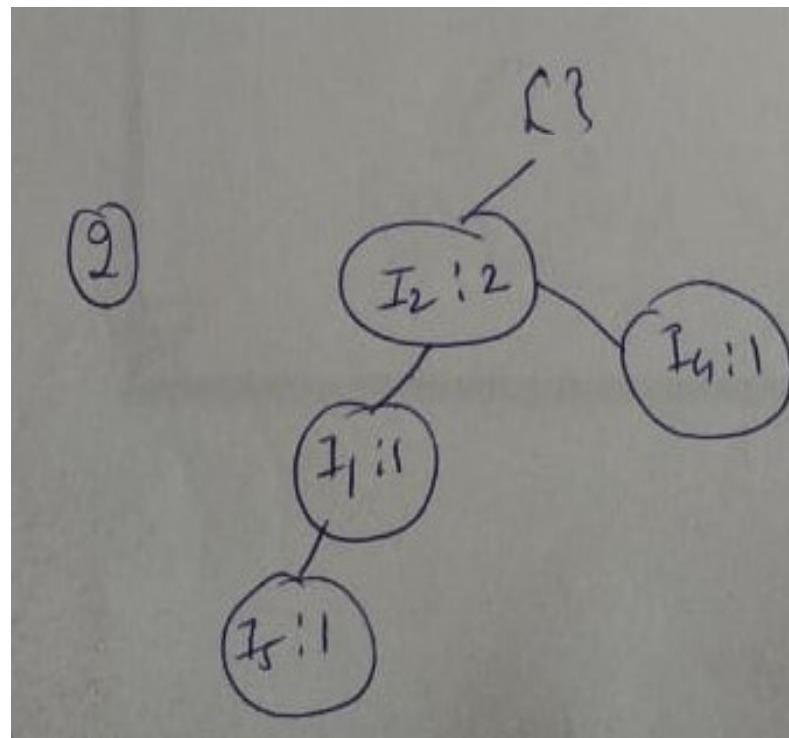
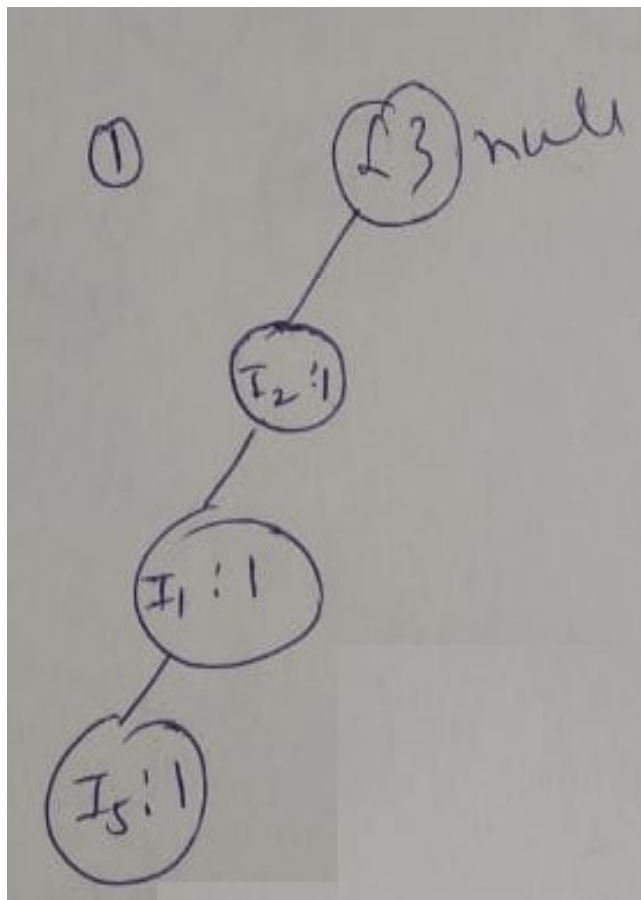
Sort the items in each transactions according to its support count

I_2 (7)
 I_1 (6)
 I_3 (6)
 I_4 (2)
 I_5^- (4)

| | |
|---|---------------------|
| 1 | $I_2 I_1 I_5^-$ |
| 2 | $I_2 I_4$ |
| 3 | $I_2 I_3$ |
| 4 | $I_2 I_1 I_4$ |
| 5 | $I_1 I_3$ |
| 6 | $I_2 I_3$ |
| 7 | $I_1 I_3$ |
| 8 | $I_2 I_1 I_3 I_5^-$ |
| 9 | $I_2 I_1 I_3$ |

Example : FP Growth cont...

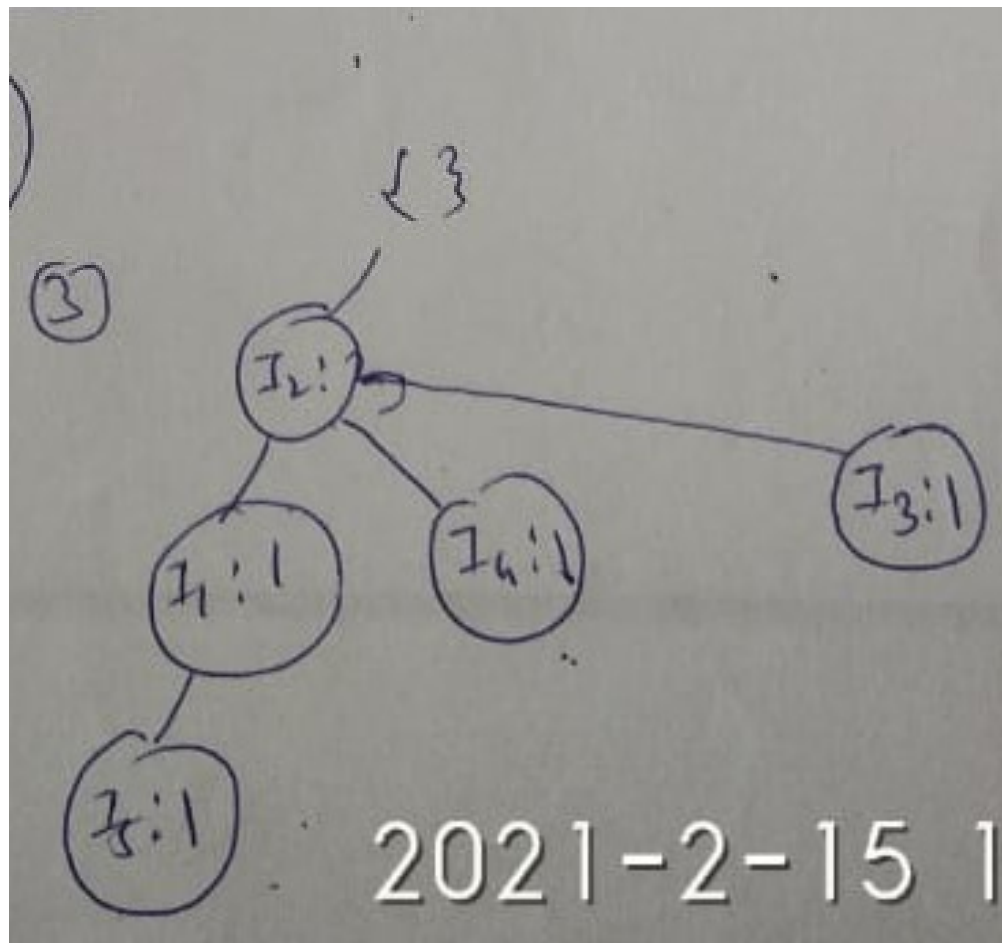
| | |
|---|---------------|
| 1 | $I_2 I_1 I_5$ |
| 2 | $I_2 I_4$ |
| 3 | $I_2 I_3$ |
| 4 | $I_2 I_1 I_4$ |



Example : FP Growth cont...

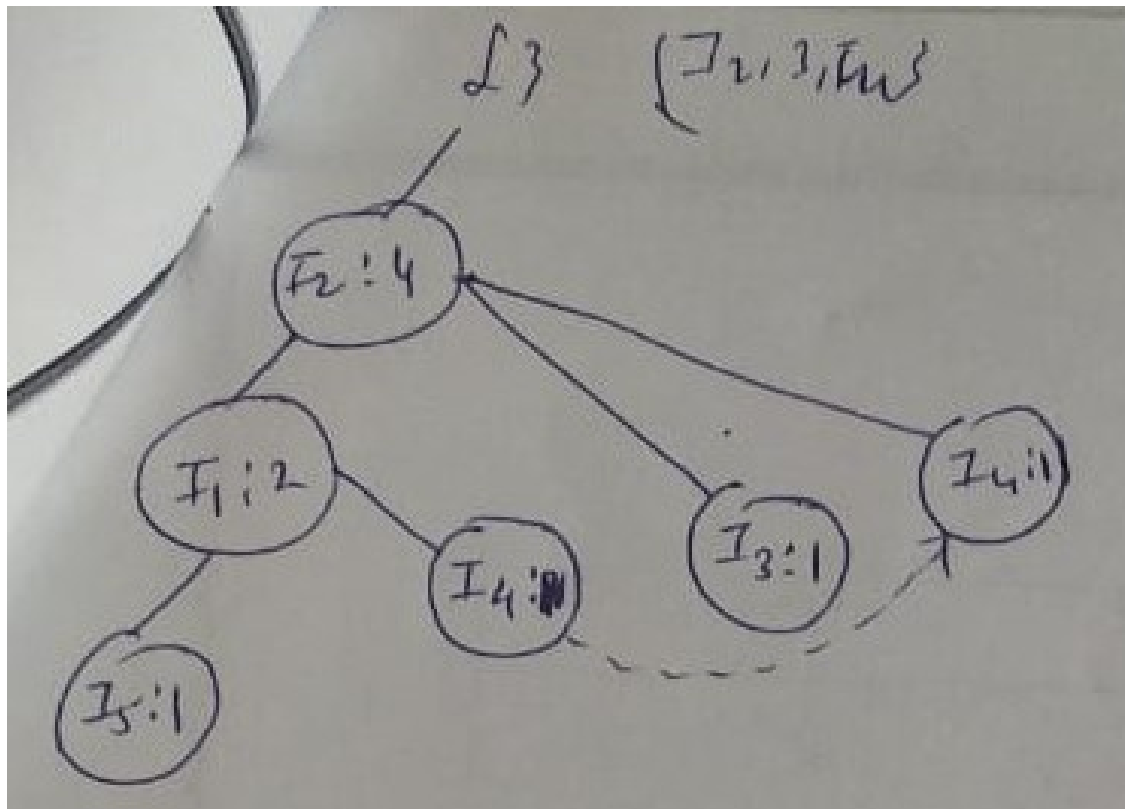
| | |
|---|-------------|
| 3 | I_2 I_3 |
|---|-------------|

| | |
|---|-------------------|
| 1 | I_2 I_1 I_5 |
| 2 | I_2 I_4 |
| 3 | I_2 I_3 |
| 4 | I_2 I_1 I_4 |

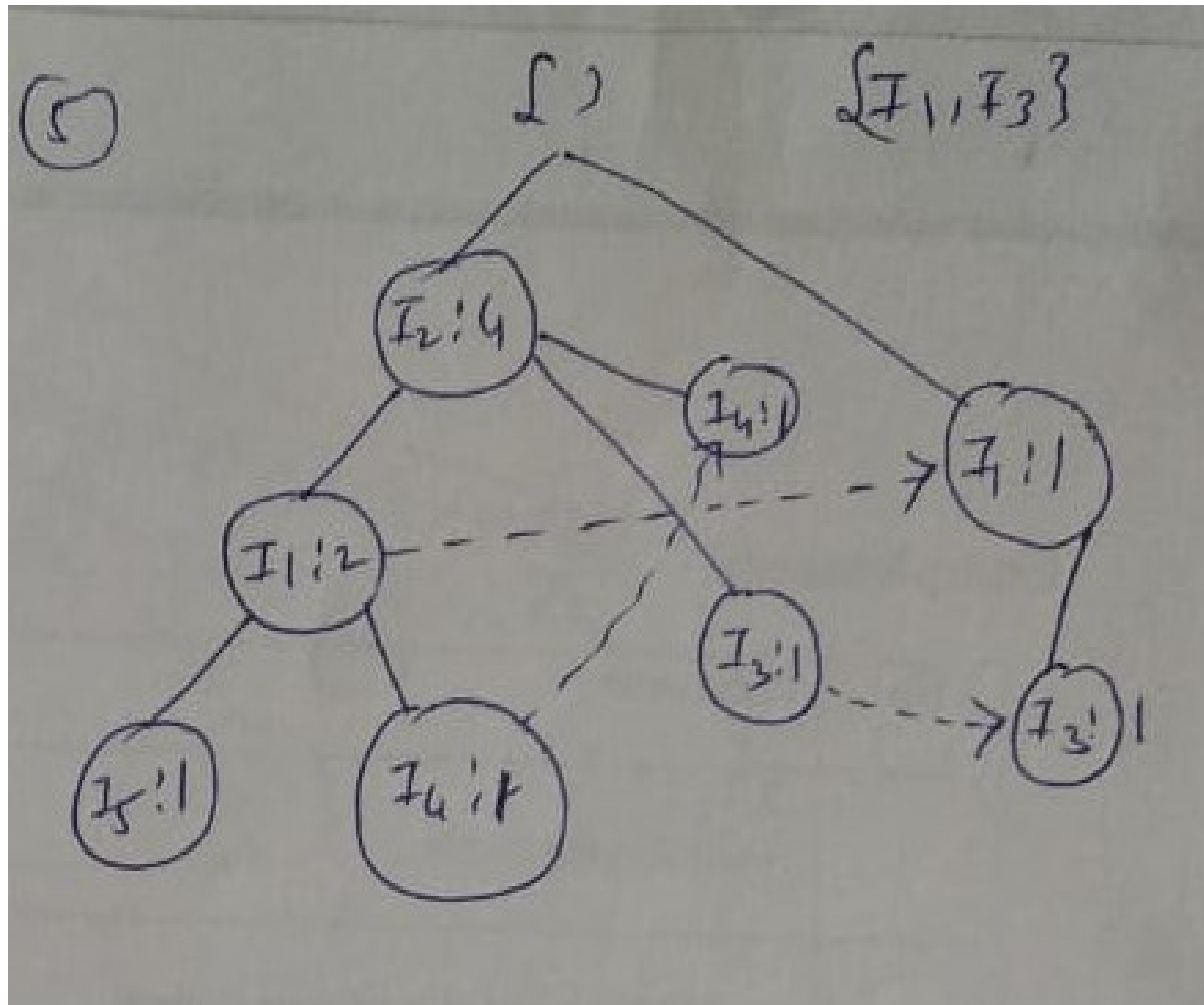


Example : FP Growth cont...

| | |
|---|---------------|
| 1 | $I_2 I_1 I_5$ |
| 2 | $I_2 I_4$ |
| 3 | $I_2 I_3$ |
| 4 | $I_2 I_1 I_4$ |



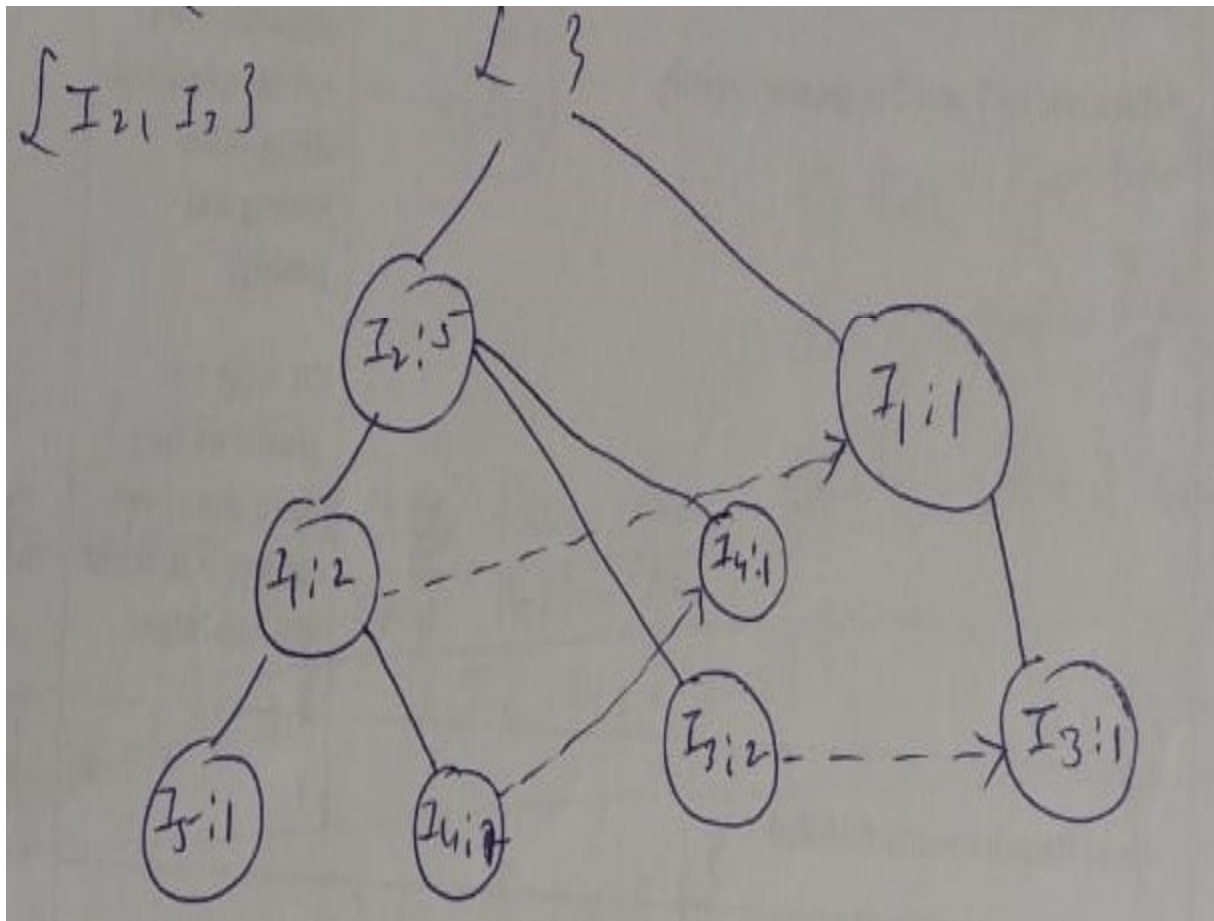
Example : FP Growth cont...



| | |
|---|-----------------|
| 1 | $I_2 I_1 I_5^-$ |
| 2 | $I_2 I_4$ |
| 3 | $I_2 I_3$ |
| 4 | $I_2 I_1 I_4$ |

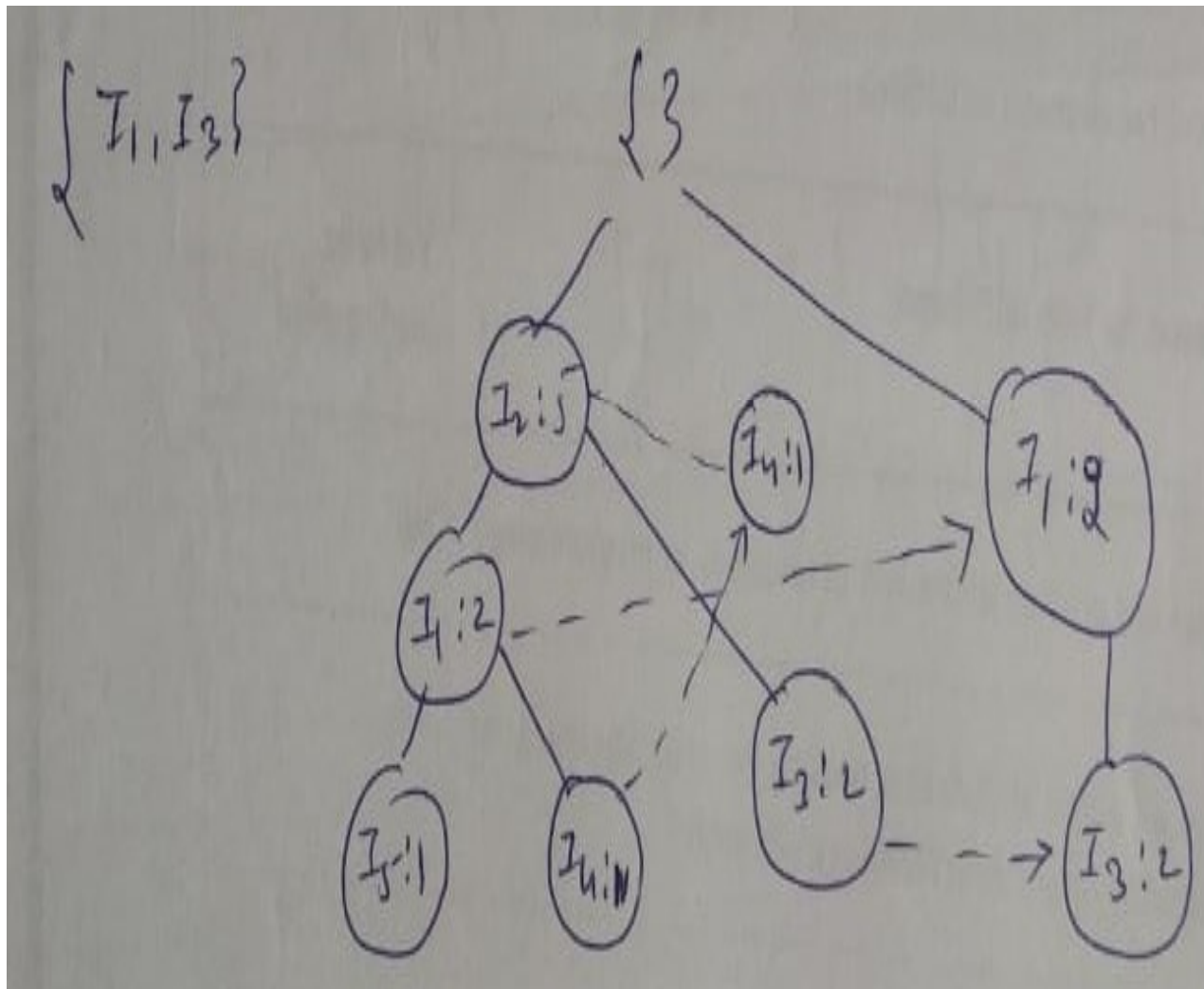
| | |
|---|---------------------|
| 5 | $I_1 I_3$ |
| 6 | $I_2 I_3$ |
| 7 | $I_1 I_3$ |
| 8 | $I_2 I_1 I_3 I_5^-$ |
| 9 | $I_2 I_1 I_3$ |

Example : FP Growth cont...



| | |
|---|-------------------|
| 5 | $I_1 I_3$ |
| 6 | $I_2 I_3$ |
| 7 | $I_1 I_3$ |
| 8 | $I_2 I_1 I_3 I_5$ |
| 9 | $I_2 I_1 I_3$ |

Example : FP Growth cont...



| | |
|---|-------------------|
| 5 | $I_1 I_3$ |
| 6 | $I_2 I_3$ |
| 7 | $I_1 I_3$ |
| 8 | $I_2 I_1 I_3 I_5$ |
| 9 | $I_2 I_1 I_3$ |

Add 8 and 9 transaction

Final FP-Tree

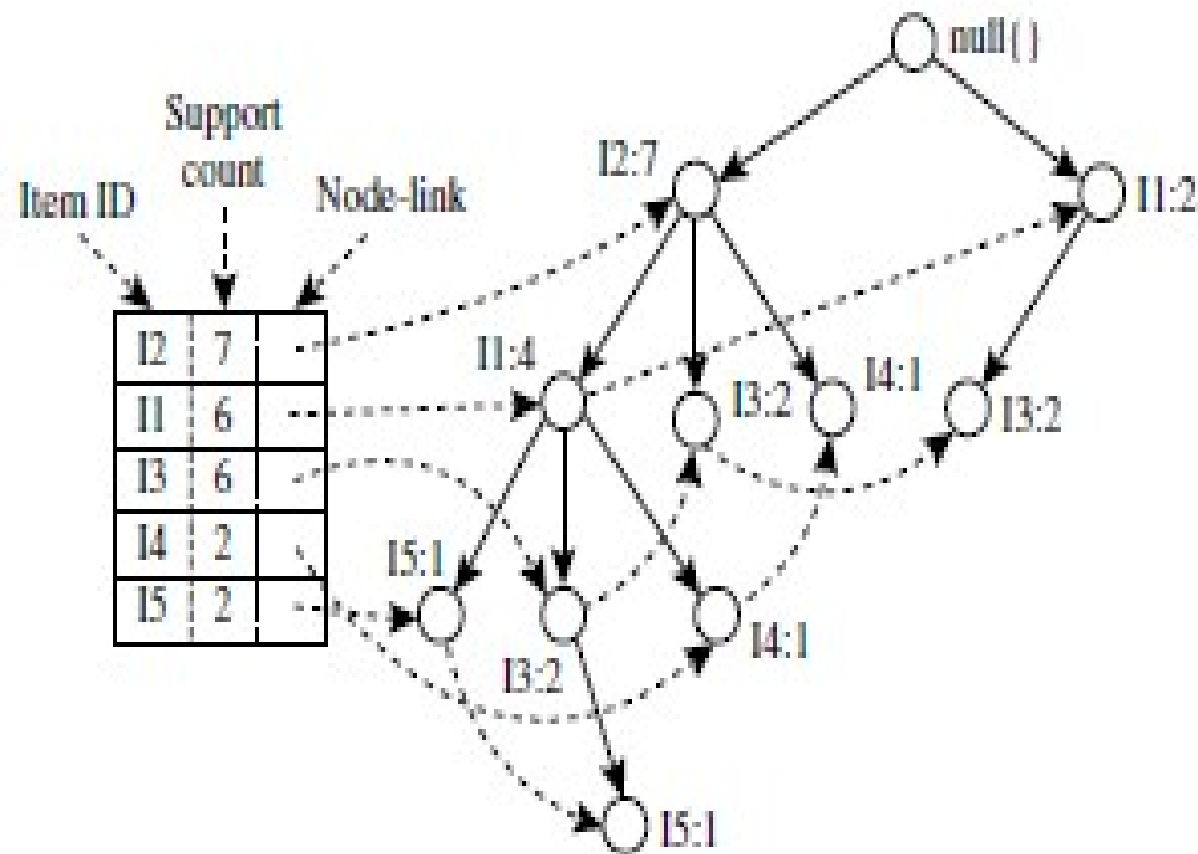
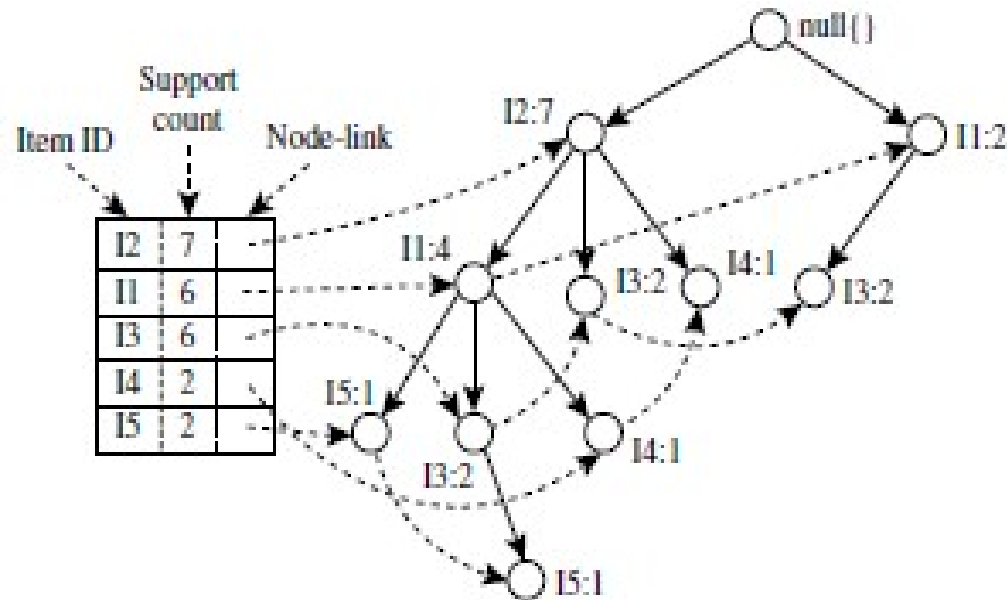


Image Source: [2]

| | |
|---|---------------------|
| 1 | $I_2 I_1 I_5^-$ |
| 2 | $I_2 I_4$ |
| 3 | $I_2 I_3$ |
| 4 | $I_2 I_1 I_4$ |
| 5 | $I_1 I_3$ |
| 6 | $I_2 I_3$ |
| 7 | $I_1 I_3$ |
| 8 | $I_2 I_1 I_3 I_5^-$ |
| 9 | $I_2 I_1 I_3$ |

Mine FP pattern from Tree



| Item | Conditional Pattern Base | Conditional FP-tree | Frequent Patterns Generated |
|------|---------------------------------|-------------------------|-------------------------------------------|
| 15 | {(12, 11: 1), (12, 11, 13: 1)} | (12: 2, 11: 2) | {12, 15: 2}, {11, 15: 2}, {12, 11, 15: 2} |
| 14 | {(12, 11: 1), (12: 1)} | (12: 2) | {12, 14: 2} |
| 13 | {(12, 11: 2), (12: 2), (11: 2)} | (12: 4, 11: 2), (11: 2) | {12, 13: 4}, {11, 13: 4}, {12, 11, 13: 2} |
| 11 | {(12: 4)} | (12: 4) | {12, 11: 4} |

Image Source: [2]

Evaluation of Association Patterns

Well accepted criteria for evaluating the quality of patterns:

- **Objective interestingness measure:**
 - Uses statistics derived from data to find interestingness.
 - E.g. **Support, confidence, lift**, and others.
- **Subjective interestingness measure:**
 - Visualization (human user in loop).
 - Template based approach (it allows the user to restrict the type of patterns extracted by algorithm using user-specified template)
 - Use domain information such as concept hierarchy(lower level of the hierarchy are considered redundant → eliminated

Criticism to Support and Confidence

- Example 1: (Aggarwal & Yu, PODS98)

- Among 5000 students
 - 3000 play basketball
 - 3750 eat cereal
 - 2000 both play basketball and eat cereal

| | basketball | not basketball | sum(row) | |
|------------|------------|----------------|----------|-----|
| cereal | 2000 | 1750 | 3750 | 75% |
| not cereal | 1000 | 250 | 1250 | 25% |
| sum(col.) | 3000 | 2000 | 5000 | |
| | 60% | 40% | | |

play basketball \Rightarrow eat cereal [40%, 66.7%]

misleading because the overall percentage of students eating cereal is 75% which is higher than 66.7%.

play basketball \Rightarrow not eat cereal [20%, 33.3%]

is more accurate, although with lower support and confidence

Objective interestingness measure

➤ Lift ($A \rightarrow B$) = confidence ($A \rightarrow B$) / support (B)

f_{1+} = support count of A

f_{+1} = support count of B

$$\text{Lift}(A \rightarrow B) = \frac{\text{sup}(A, B)}{\text{sup}(A) \cdot \text{sup}(B)} = \frac{P(B | A)}{P(B)}$$

- A and B negatively correlated, if the value is less than 1; otherwise A and B positively correlated

1. If the rule had a lift of 1, it would imply that the probability of occurrence of the antecedent and that of the consequent are independent of each other. Hence, no rule can be drawn involving those two events.
2. If the lift is > 1 , that lets us know the degree to which those two occurrences are dependent on one another, and makes those rules potentially useful for predicting the consequent in future data sets.
3. If the lift is < 1 , that lets us know the items are substitute to each other. This means that presence of one item has negative effect on presence of other item and vice versa.

Image Source: [3,5]

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

Beverage preferences among a group of 1000 people.

| | Coffee | $\bar{\text{Coffee}}$ | |
|--------------------|--------|-----------------------|------|
| Tea | 150 | 50 | 200 |
| $\bar{\text{Tea}}$ | 650 | 150 | 800 |
| | 800 | 200 | 1000 |

Lift of a Rule

■ Example 1 (cont)

■ *play basketball* \Rightarrow *eat cereal* [40%, 66.7%]

$$\text{LIFT} = \frac{\frac{2000}{5000}}{\frac{3000}{5000} \times \frac{3750}{5000}} = 0.89$$

■ *play basketball* \Rightarrow *not eat cereal* [20%, 33.3%]

$$\text{LIFT} = \frac{\frac{1000}{5000}}{\frac{3000}{5000} \times \frac{1250}{5000}} = 1.33$$

| | basketball | not basketball | sum(row) |
|------------|------------|----------------|----------|
| cereal | 2000 | 1750 | 3750 |
| not cereal | 1000 | 250 | 1250 |
| sum(col.) | 3000 | 2000 | 5000 |

Image Source: [5]

Problems With Lift

- Rules that hold 100% of the time may not have the highest possible lift. For example, if 5% of people are Vietnam veterans and 90% of the people are more than 5 years old, we get a lift of $0.05/(0.05*0.9)=1.11$ which is only slightly above 1 for the rule
- Vietnam veterans \rightarrow more than 5 years old.
- And, lift is *symmetric*:
- not eat cereal* \Rightarrow *play basketball* [20%, 80%]

$$\text{LIFT} = \frac{\frac{1000}{5000}}{\frac{1250}{5000} * \frac{3000}{5000}} = 1.33$$

Objective interestingness measure

➤ Conviction = $(f_{1+} f_{+0}) / (N f_{10})$

$$\text{conviction}(A \rightarrow C) = \frac{1 - \text{support}(C)}{1 - \text{confidence}(A \rightarrow C)}$$

It is the ratio of the expected frequency that X occurs without Y (that is to say, the frequency that the rule makes an incorrect prediction) if X and Y were independent divided by the observed frequency of incorrect predictions.

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

Beverage preferences among a group of 1000 people.

| | Coffee | $\overline{\text{Coffee}}$ | |
|-------------------------|--------|----------------------------|------|
| Tea | 150 | 50 | 200 |
| $\overline{\text{Tea}}$ | 650 | 150 | 800 |
| | 800 | 200 | 1000 |

Image Source: [3]

$$\text{Conv}(A \rightarrow B) = \frac{\text{sup}(A) \cdot \text{sup}(\bar{B})}{\text{sup}(A, \bar{B})} = \frac{P(A) \cdot P(\bar{B})}{P(A, \bar{B})} = \frac{P(A)(1 - P(B))}{P(A) - P(A, B)}$$

- Conviction is a measure of the implication and has value 1 if items are unrelated.

- play basketball* \Rightarrow *eat cereal* [40%, 66.7%]

- eat cereal* \Rightarrow *play basketball* conv:0.85

$$\text{Conv} = \frac{\frac{3000}{5000} \cdot 1 - \frac{3750}{5000}}{\frac{3000}{5000} - \frac{2000}{5000}} = 0.75$$

- play basketball* \Rightarrow *not eat cereal* [20%, 33.3%]

- not eat cereal* \Rightarrow *play basketball* conv:1.43

$$\text{Conv} = \frac{\frac{3000}{5000} \cdot 1 - \frac{1250}{5000}}{\frac{3000}{5000} - \frac{1000}{5000}} = 1.125$$

Summary of FP-Growth Algorithm

- Mining frequent patterns can be viewed as first mining 1-itemset and progressively growing each 1-itemset by mining on its **conditional pattern base** recursively
 - Transform a frequent k-itemset mining problem into a sequence of k frequent 1-itemset mining problems via a set of conditional pattern bases
-

References

- [1] Han, Jiawei, Jian Pei, and Yiwen Yin. "Mining frequent patterns without candidate generation." *ACM sigmod record* 29.2 (2000): 1-12.
- [2] Han, Jiawei, Jian Pei, and Micheline Kamber. *Data mining: concepts and techniques*. Elsevier, 2012(Third Edition).
- [3] Tan, Pang-Ning, Michael Steinbach, and Vipin Kumar. Introduction to data mining. Pearson Education India, 2016.
- [4] .http://www.cs.sunysb.edu/~cse634/lecture_notes/07apriori.pdf
- [5] https://paginas.fe.up.pt/~ec/files_0506/slides/04_AssociationRules.pdf