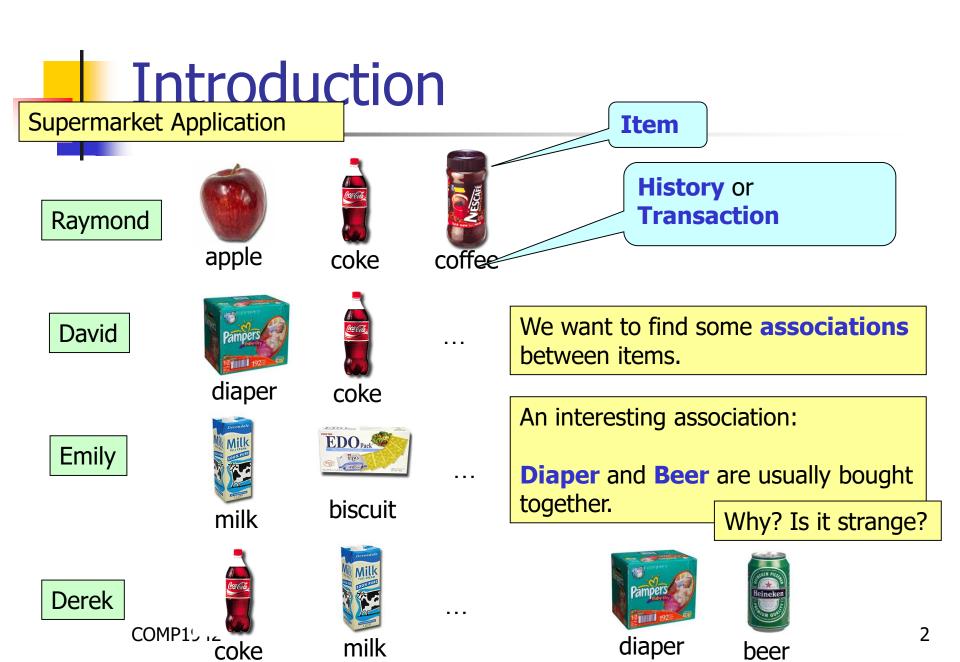


Association Rule Mining (Concepts)

Prepared by Raymond Wong Presented by Raymond Wong raywong@cse





An interesting association:

Diaper and **Beer** are usually bought together.

Why? Is it strange?





COMP1942



An interesting association:

Diaper and **Beer** are usually bought together.

Why? Is it strange?





diaper

beer

Reasons:

This pattern occurs frequently in the **early evening**.





An interesting association:

Diaper and **Beer** are usually bought together.

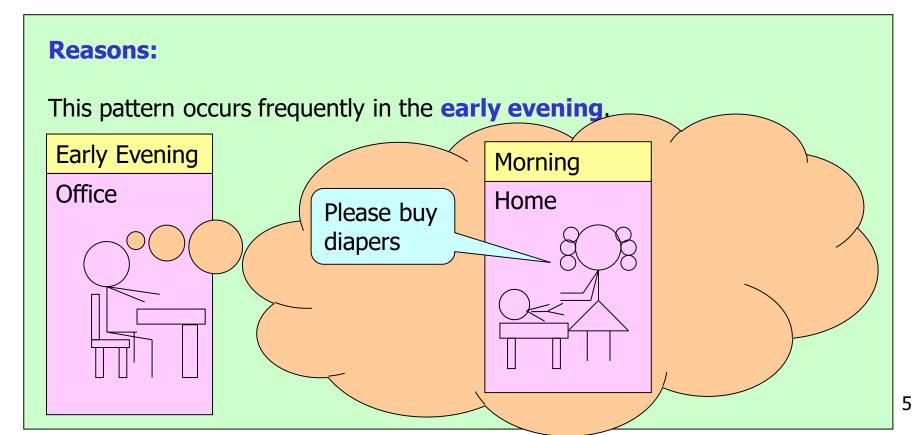
Why? Is it strange?





diaper

beer



Introduction

- Applications of Association Rule Mining
 - Supermarket
 - Web Mining
 - Medical analysis
 - Bioinformatics
 - Network analysis (e.g., Denial-of-service (DoS))
 - Programming Pattern Finding

Outline

- Association Rule Mining
 - Problem Definition
- Algorithm
- How to use the data mining tool

| TID | Α | В | С | D | E |
|-----|---|---|---|---|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

В

С

D

Ε

Itemsets:

{B, C}

{A, B, C}

{B, C, D}

{A}

COMP1 2-itemset

3-itemset

3-itemset

1-itemset

Large itemsets:

itemsets with support >= a threshold (e.g., 3)

Frequent itemsets

on Ruis

Mining

| TID | Α | В | С | D | Е |
|-----|---|---|---|---|----|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1/ |

e.g., {A}, {B}, {B, C} but NOT {A, B, C}

Support = 3

Support = 4

Single Items (or simply items):

В

С

D

Ε

Itemsets:

{B, C}

 $\{A, B, C\}$

{B, C, D}

{A}

1-frequent itemset of size 3

COMP1 Support = 3

Support = 1

3-frequent itemset of size 2

| TID | Α | В | С | D | E |
|-----|---|-----|-----|-----|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1_ | 1 | 0 | 0 |
| t4 | 1 | (1_ | _ 1 | 1 (| 1 |
| t5 | 0 | (1 | 1 | 0 (| 1 |

Support = 2

Association rules:

$$\{B, C\} \rightarrow E$$

| TID | Α | В | U | D | Е |
|-----|---|-----|-----|---|----|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | (1_ | _ 1 | 0 | 0 |
| t4 | 1 | (1 | 1 | 1 | 1 |
| t5 | 0 | (1 | 1 | 0 | 1/ |

Support = 2

Confidence = 2/3 = 66.7%

Association rules:

 $\{B, C\} \rightarrow E$

| TID | Α | В | C | D | Е |
|-----|---|-----|-----|---|----|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | (1_ | _ 1 | 0 | 0 |
| t4 | 1 | (1 | 1 | 1 | 1 |
| t5 | 0 | (1 | 1 | 0 | 1/ |
| | | | | | |

Support = 2

Confidence = 2/3 = 66.7%

Association rules:

$$\{B, C\} \rightarrow E$$

Support = 3

$$B \rightarrow C$$

| TID | Α | В | С | D | Е |
|-----|---|-----|---|---|----|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | (1) | 0 | 1 | 1 |
| t3 | 0 | (1) | 1 | 0 | 0 |
| t4 | 1 | | 1 | 1 | 1 |
| t5 | 0 | 1, | 1 | 0 | 1/ |

Support = 2

Confidence = 2/3 = 66.7%

Association rules:

$$\{B, C\} \rightarrow E$$

Support = 3

 $B \rightarrow C$

Confidence= 3/4 = 75%

Antecedent

Consequent

Association rules with

- 1. Support >= a threshold (e.g., 3)
- 2. Confidence >= another threshold (e.g., 50%)

| ASSOCI | atı |
|--------|-----|
| | |

| TID | А | В | С | D | Е |
|-----|---|---|---|---|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

How can we find all "interesting" association rules?

Step 2: to find all "interesting" rules after Step 1

from all "large" itemsets
find the association rule with confidence
= 50%

Problem:

We want to find some "interesting" association rules

$$\{B, C\} \rightarrow E$$
Support = 2

Confidence = 2/3 = 66.7%

$$B \rightarrow C$$

Support
$$= 3$$

Confidence =
$$3/4 = 75\%$$

- The previous definition is the traditional definition of association rule mining
- The following is a new concept used in the data mining tool
 - Lift Ratio

Confidence of the rule

Lift Ratio =

Expected Confidence of the consequent of the rule

| TID | Α | В | С | D | Е |
|-----|---|---|---|---|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

Confidence of the rule = 2/3 = 66.7%

Association rules:

$$\{B, C\} \rightarrow E$$

Expected Confidence of the consequent of the rule = 3/5 = 60%

Lift Ratio of the rule = 66.7/60 = 1.11

| TID | Α | В | С | D | Е |
|-----|---|---|---|---|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

Confidence of the rule = 3/4 = 75%

Association rules:

$$B \rightarrow C$$

Expected Confidence of the consequent of the rule = 3/5 = 60%

Lift Ratio of the rule = 75/60 = 1.25

Outline

- Association Rule Mining
 - Problem Definition
- Algorithm
- How to use the data mining tool

Algorithm

- Two major algorithms
 - Apriori Approach
 - FP-growth Approach



Apriori Algorithm

- Properties
- Algorithm

Suppose we want to find all "large" itemsets (e.g., itemsets with support >= 3)

| TID | Α | В | С | D | Е |
|-----|---|---|---|---|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

{B, C} is large

Support of $\{B, C\} = 3$

Is {B} large?

Is {C} large?

Property 1: If an itemset S is large, then any proper subset of S must be large.



Suppose we want to find all "large" itemsets (e.g., itemsets with support >= 3)

| TID | Α | В | С | D | Е |
|-----|---|---|---|---|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

{B, C, E} is NOT large

Support of $\{B, C, E\} = 2$

Is {A, B, C, E} large?

Is {B, C, D, E} large?

Property 2: If an itemset S is NOT large, then any proper superset of S must NOT be large.

Property 1: If an itemset S is large, then any proper subset of S must be large.

Property 2: If an itemset S is NOT large, then any proper superset of S must NOT be large.

COMP1942 23



Apriori Algorithm

- Properties
- Algorithm

| TID | Α | В | С | D | Е |
|-----|---|---|---|---|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

| Item | Count |
|------|-------|
| Α | 3 |
| В | |
| С | |
| D | |
| Е | |

COMP1942 25

Suppose we want to find all "large" itemsets (e.g., itemsets with support >= 3)

| TID | Α | В | С | D | Е |
|-----|---|---|---|---|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

| Item | Count | | |
|------|-------|--|--|
| А | 3 | | |
| В | 4 | | |
| С | 3 | | |
| D | 3 | | |
| Е | (3) | | |

Thus, {A}, {B}, {C}, {D} and {E} are "large" itemsets of size 1 (or, "large" 1-itemsets).

We set $L_1 = \{\{A\}, \{B\}, \{C\}, \{D\}, \{E\}\}\}$

Suppose we want to find all "large" itemsets (e.g., itemsets with support >= 3)

| TID | A | В | U | D | Large 2-i |
|-----|---|---|---|---|-----------|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

Candidate Generation

C₂

 L_1

 L_2

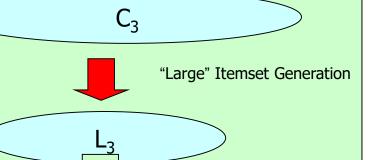
"Large" Itemset Generation

Large 3-itemset Generation

Candidate Generation

Thus, {A}, {B}, {C}, {D} and {E} are "large" itemsets of size 1 (or, "large" 1-itemsets).

We set $L_1 = \{\{A\}, \{B\}, \{C\}, \{D\}, \{E\}\}\}$





Suppose we want to find all "large" itemse 1. itemsets with support >= 3)

- 1. Join Step
- 2. Prune Step

| TID | A | В | С | D | Large 2-itemset Generatio | |
|-----|---|---|---|---|---------------------------|--|
| t1 | 1 | 0 | 0 | 1 | 0 | |
| t2 | 1 | 1 | 0 | 1 | 1 | |
| t3 | 0 | 1 | 1 | 0 | 0 | |
| t4 | 1 | 1 | 1 | 1 | 1 | |
| t5 | 0 | 1 | 1 | 0 | 1 | |

Candidate Generation

 C_2

 L_1

"Large" Itemset Generation

Large 3-itemset Generation

Counting Step

Candidate Generation

Thus, {A}, {B}, {C}, {D} and {E} are "large" itemsets of size 1 (or, "large" 1-itemsets).

We set $L_1 = \{\{A\}, \{B\}, \{C\}, \{D\}, \{E\}\}\}$

 L_2

"Large" Itemset Generation

L₃



Candidate Generation

- Join Step
- Prune Step

Property 1: If an itemset S is large, then any proper subset of S must be large.

Property 2: If an itemset S is NOT large, then any proper superset of S must NOT be large.

Join Step

| TID | Α | В | С | D | Е |
|-----|---|---|---|---|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

Suppose we know that itemset $\{B, C\}$ and itemset $\{B, E\}$ are large (i.e., L_2).

It is possible that itemset $\{B, C, E\}$ is also large (i.e., C_3).

Join Step

- Join Step
 - Input: L_{k-1}, a set of all large (k-1)-itemsets
 - Output: C_k, a set of candidates k-itemsets
 - Algorithm:

```
• insert into C<sub>k</sub> select p.item<sub>1</sub>, p.item<sub>2</sub>, ..., p.item<sub>k-1</sub>, q.item<sub>k-1</sub> from L<sub>k-1</sub> p, L<sub>k-1</sub> q where p.item<sub>1</sub> = q.item<sub>1</sub>, p.item<sub>2</sub> = q.item<sub>2</sub>, ... p.item<sub>k-2</sub> = q.item<sub>k-2</sub>, p.item<sub>k-1</sub> < q.item<sub>k-1</sub>
```

Property 1: If an itemset S is large, then any proper subset of S must be large.

Property 2: If an itemset S is NOT large, then any Prune Step proper superset of S must NOT be large.

| TID | Α | В | С | D | Е |
|-----|---|---|---|---|---|
| t1 | 1 | 0 | 0 | 1 | 0 |
| t2 | 1 | 1 | 0 | 1 | 1 |
| t3 | 0 | 1 | 1 | 0 | 0 |
| t4 | 1 | 1 | 1 | 1 | 1 |
| t5 | 0 | 1 | 1 | 0 | 1 |

Suppose we know that itemset $\{B, C\}$ and itemset $\{B, E\}$ are large (i.e., L_2).

It is possible that itemset $\{B, C, E\}$ is also large (i.e., C_3).

Suppose we know that {C, E} is not large. We can prune $\{B, C, E\}$ in C_3 .

Prune Step

- Prune Step
 - forall itemsets c ∈C_k (from Join Step) do
 - for all (k-1)-subsets s of c do
 - if (s not in L_{k-1}) then
 - delete c from C_k



Suppose we want to find all "large" itemse itemsets with support >= 3)

- . Join Step
- 2. Prune Step

Counting Step

Candidate Generation

| | | _ | _ | | | |
|-----|---|----------|--------------------|---|----------------------------|--|
| TID | Α | В | С | D | Large 2-itemset Generation | |
| t1 | 1 | 0 | 0 | 1 | 0 | |
| t2 | 1 | 1 | 0 | 1 | 1 | |
| t3 | 0 | 1 | 1 | 0 | 0 | |
| t4 | 1 | 1 | 1 | 1 | 1 | |
| t5 | 0 | 1 | 1 | 0 | 1 | |
| | | Large 3- | itemset Generation | | | |

Candidate Generation

C2

"Large" Itemset Generation

 L_1

Thus, {A}, {B}, {C}, {D} and {E} are "large" itemsets of size 1 (or, "large" 1-itemsets).

We set $L_1 = \{\{A\}, \{B\}, \{C\}, \{D\}, \{E\}\}\}$

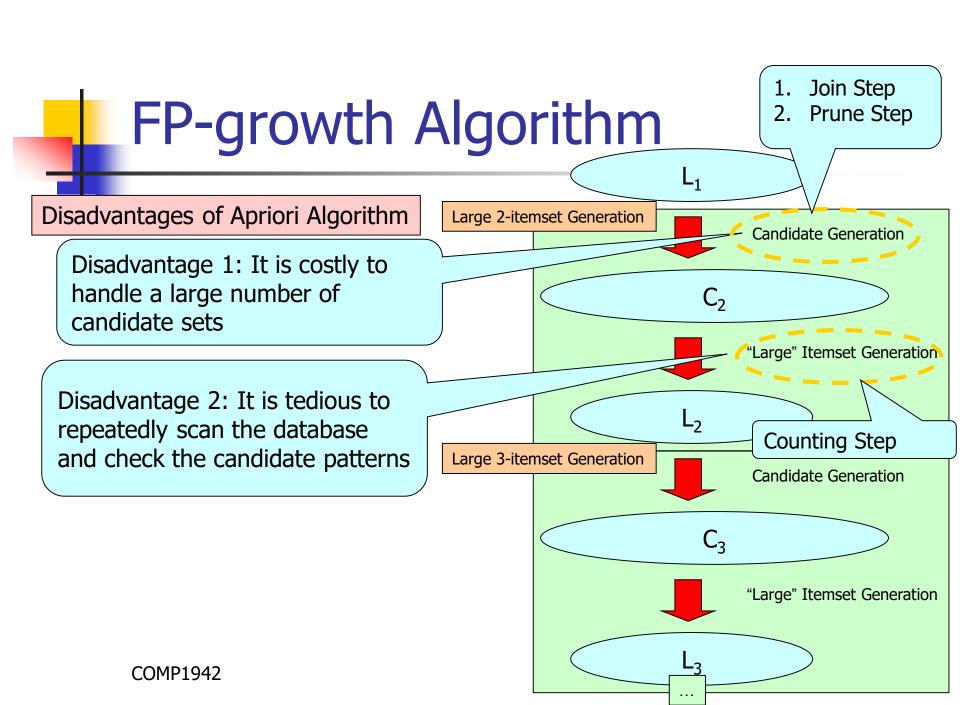
"Large" Itemset Generation

Counting Step

- After the candidate generation (i.e., Join Step and Prune Step), we are given a set of candidate itemsets
- We need to verify whether these candidate itemsets are large or not
- We have to scan the database to obtain the count of each itemset in the candidate set.
- Algorithm
 - For each itemset c in C_k,
 - obtain the count of c (from the database)
 - If the count of c is smaller than a given threshold,
 - remove it from C_k
 - The remaining itemsets in C_k correspond to L_k

Algorithm

- Two major algorithms
 - Apriori Approach
 - FP-growth Approach.





- Scan the database once to store all essential information in a data structure called FP-tree (Frequent Pattern Tree)
- The FP-tree is concise and is used in directly generating large itemsets

COMP1942

Step 1: Deduce the ordered frequent items. For items with the same frequency, the order is given by the alphabetical order.

Step 2: Construct the FP-tree from the above data

Step 3: From the FP-tree above, construct the FP-

conditional tree for each item (or itemset).

Step 4: Determine the frequent patterns.

COMP1942

Frequent Itemset Mining

Problem: to find all "large" (or frequent) itemsets with support at least a threshold (i.e., itemsets with support >= 3)

| TID | Items Bought |
|-----|------------------------|
| 100 | a, b, c, d, e, f, g, h |
| 200 | a, f, g |
| 300 | b, d, e, f, j |
| 400 | a, b, d, i, k |
| 500 | a, b, e, g |

| TID | Items Bought | |
|-----|------------------------|--|
| 100 | a, b, c, d, e, f, g, h | |
| 200 | a, f, g | |
| 300 | b, d, e, f, j | |
| 400 | a, b, d, i, k | |
| 500 | a, b, e, g | |

| TID | Items Bought | |
|-----|------------------------|--|
| 100 | a, b, c, d, e, f, g, h | |
| 200 | a, f, g | |
| 300 | b, d, e, f, j | |
| 400 | a, b, d, i, k | |
| 500 | a, b, e, g | |

COMP1942

42

| TID | Items Bought | (Ordered) Frequent Items |
|-------|------------------------|--------------------------|
| 100 (| a, b, c, d, e, f, g, h | |
| 200 (| a, f, g | |
| 300 | b, d, e, f, j | |
| 400 (| a, b, d, i, k | |
| 500 (| a, b, e, g | |

| 300 a, 5, c, g | | |
|----------------|-----------|--|
| Item | Frequency | |
| а | 4 | |
| b | | |
| С | | |
| d | | |
| е | | |
| f | | |
| g | | |
| h | | |
| i | | |
| j | | |
| k | | |

| TID | Items Bought | (Ordered) Frequent Items |
|-------|------------------------|--------------------------|
| 100 | a, b, c, d, e, f, g, h | |
| 200 | a, f, g | |
| 300 (| b, d, e, f, j | |
| 400 | a, b, d, i, k | |
| 500 | a, b, e, g | |

| 3 3 3 3 | | |
|---------|-----------|--|
| Item | Frequency | |
| а | 4 | |
| b | 4 | |
| С | 1 | |
| d | 3 | |
| е | 3 | |
| f | 3 | |
| g | 3 | |
| h | 1 | |
| i | 1 | |
| j | 1 | |
| k | 1 | |

| TID | Items Bought | (Ordered) Frequent Items |
|-----|------------------------|--------------------------|
| 100 | a, b, c, d, e, f, g, h | |
| 200 | a, f, g | |
| 300 | b, d, e, f, j | |
| 400 | a, b, d, i, k | |
| 500 | a, b, e, g | |

| Threshol | d = 3 |
|----------|-------|
| | |

| Item | Frequency |
|------|-----------|
| а | 4 |
| b | 4 |
| С | 1 |
| d | (3_) |
| е | (3=) |
| f | 3_ |
| g | 3 |
| h | 1 |
| i | 1 |
| j | 1 |
| k | 1 |

| Item | Frequency |
|------|-----------|
| а | 4 |
| b | 4 |
| d | 3 |
| е | 3 |
| f | 3 |
| g | 3 |

| TID | Items Bought | (Ordered) Frequent Items |
|-----|------------------------|--------------------------|
| 100 | a, b, c, d, e, f, g, h | a, b, d, e, f, g |
| 200 | a, f, g | a, f, g |
| 300 | b, d, e, f, j | b, d, e, f |
| 400 | a, b, d, i, k | a, b, d |
| 500 | a, b, e, g | a, b, e, g |

| Item | Frequency |
|------|-----------|
| а | 4 |
| b | 4 |
| С | 1 |
| d | 3 |
| е | 3 |
| f | 3 |
| g | 3 |
| h | 1 |
| i | 1 |
| j | 1 |
| k | 1 |

| Item | Frequency |
|------|-----------|
| а | 4 |
| b | 4 |
| d | 3 |
| е | 3 |
| f | 3 |
| g | 3 |

Step 1: Deduce the ordered frequent items. For items with the same frequency, the order is given by the alphabetical order.

Step 2: Construct the FP-tree from the above data

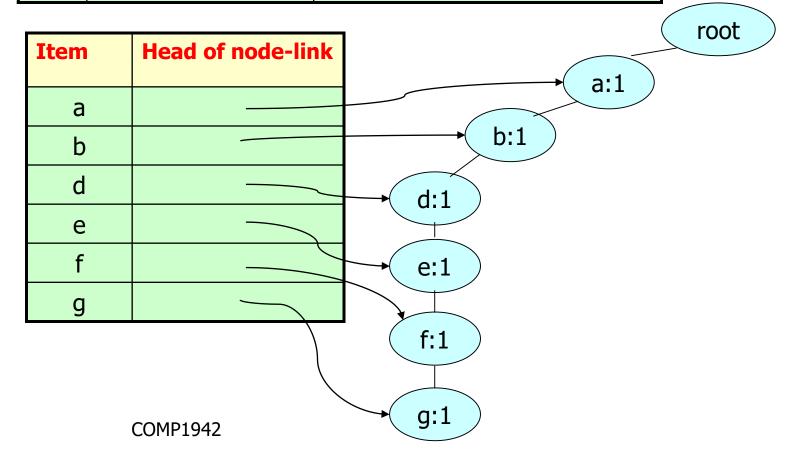
Step 3: From the FP-tree above, construct the FP-

conditional tree for each item (or itemset).

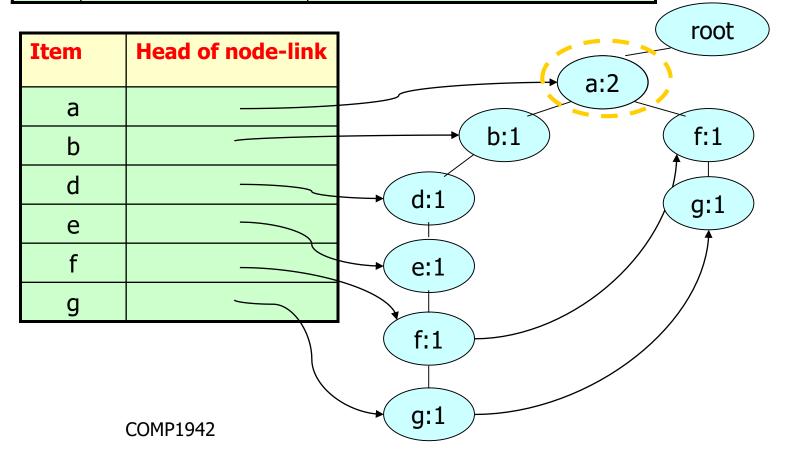
Step 4: Determine the frequent patterns.

COMP1942

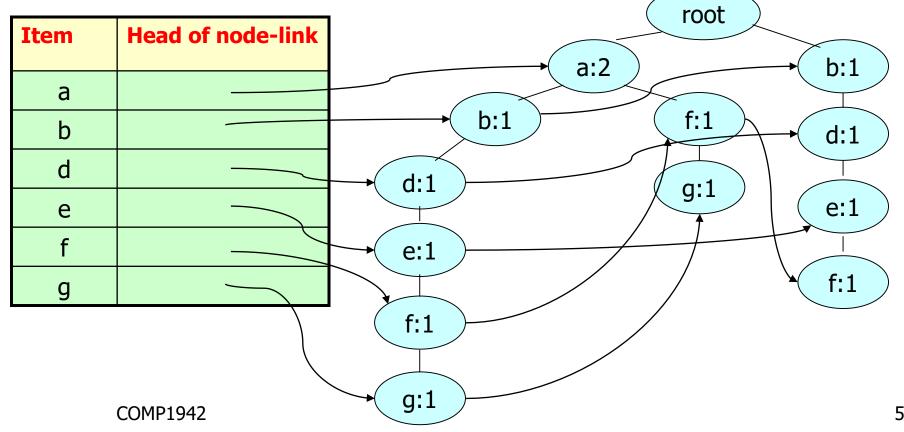
(Ordered) Frequent Items TID Items Bought a, b, d, e, f, g 100 a, b, c, d, e, f, g, h a, f, g 200 a, f, g b, d, e, f 300 b, d, e, f, j a, b, d 400 a, b, d, i, k 500 a, b, e, g a, b, e, g



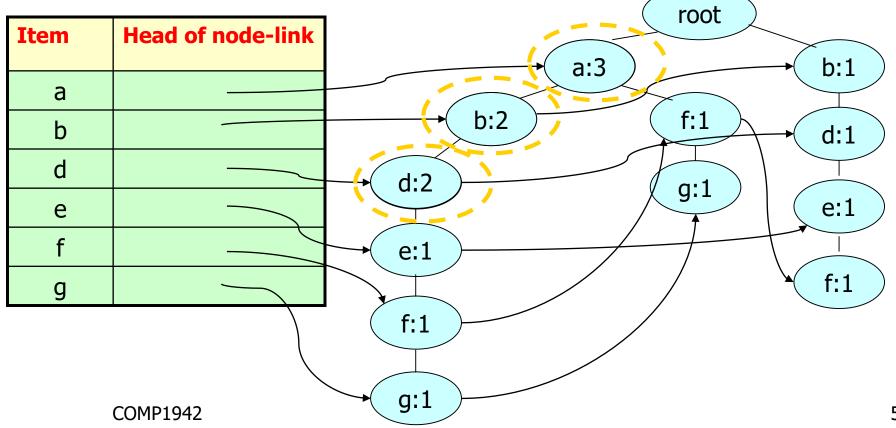
(Ordered) Frequent Items **TID Items Bought** 100 a, b, c, d, e, f, g, h a, b, d, e, f, g a, f, g 200 a, f, g b, d, e, f 300 b, d, e, f, j a, b, d 400 a, b, d, i, k 500 a, b, e, g a, b, e, g



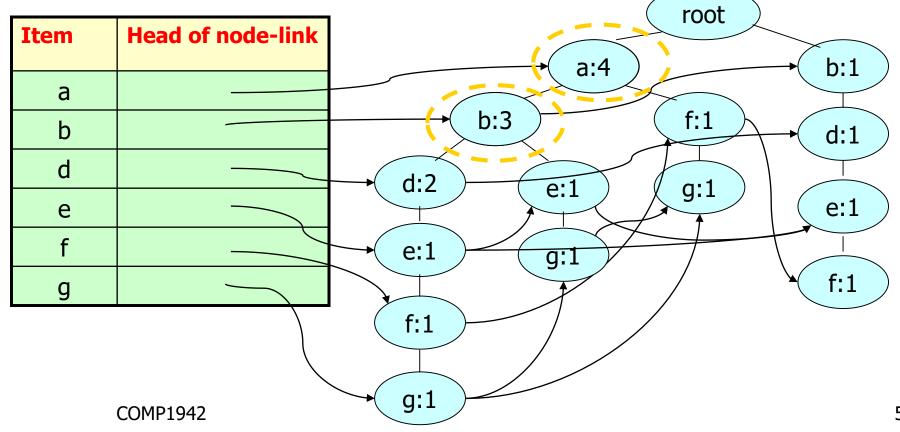
| TID | Items Bought | (Ordered) Frequent Items |
|-----|------------------------|--------------------------|
| 100 | a, b, c, d, e, f, g, h | a, b, d, e, f, g |
| 200 | a, f, g | a, f, g |
| 300 | b, d, e, f, j | b, d, e f |
| 400 | a, b, d, i, k | a, b, d |
| 500 | a, b, e, g | a, b, e, g |



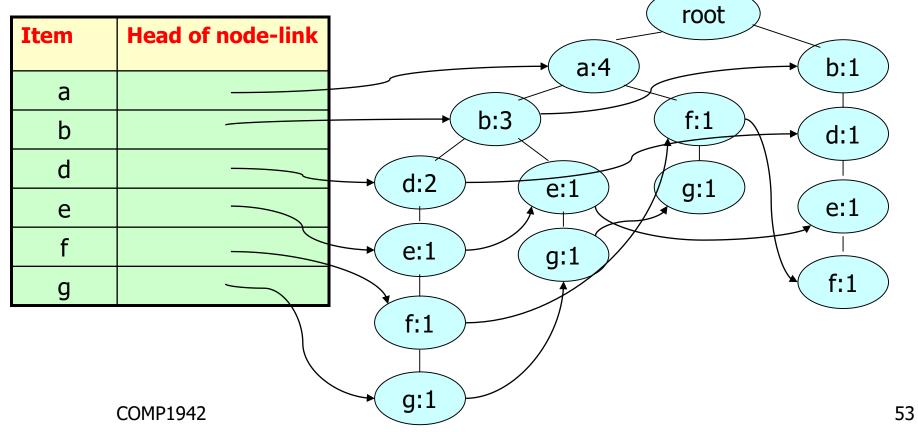
| TID | Items Bought | (Ordered) Frequent Items |
|-----|------------------------|--------------------------|
| 100 | a, b, c, d, e, f, g, h | a, b, d, e, f, g |
| 200 | a, f, g | a, f, g |
| 300 | b, d, e, f, j | b, d, e, f |
| 400 | a, b, d, i, k | a, b, d |
| 500 | a, b, e, g | a, b, e, g |



| TID | Items Bought | (Ordered) Frequent Items |
|-----|------------------------|--------------------------|
| 100 | a, b, c, d, e, f, g, h | a, b, d, e, f, g |
| 200 | a, f, g | a, f, g |
| 300 | b, d, e, f, j | b, d, e, f |
| 400 | a, b, d, i, k | a, b, d |
| 500 | a, b, e, g | a, b, e, g |



| TID | Items Bought | (Ordered) Frequent Items |
|-----|------------------------|--------------------------|
| 100 | a, b, c, d, e, f, g, h | a, b, d, e, f, g |
| 200 | a, f, g | a, f, g |
| 300 | b, d, e, f, j | b, d, e, f |
| 400 | a, b, d, i, k | a, b, d |
| 500 | a, b, e, g | a, b, e, g |



Step 1: Deduce the ordered frequent items. For items with the same frequency, the order is given by the alphabetical order.

Step 2: Construct the FP-tree from the above data

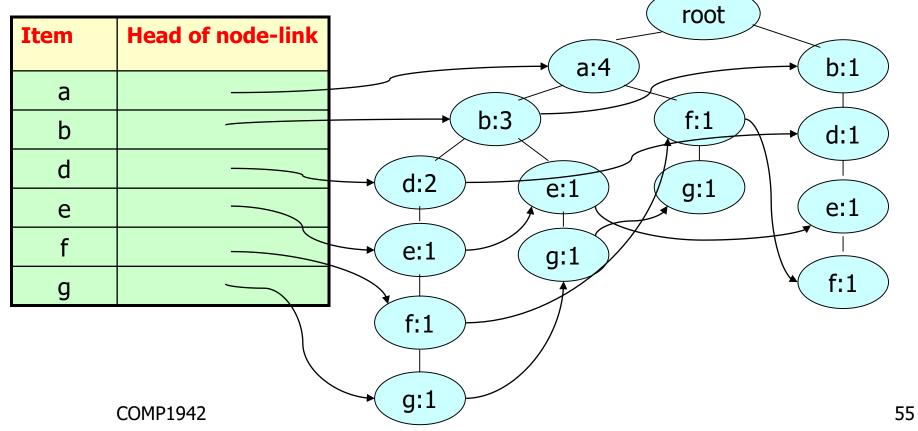
Step 3: From the FP-tree above, construct the FP-

conditional tree for each item (or itemset).

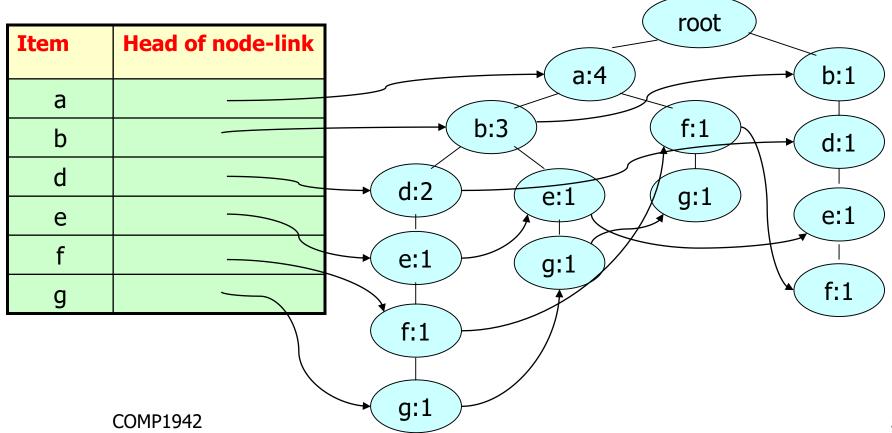
Step 4: Determine the frequent patterns.

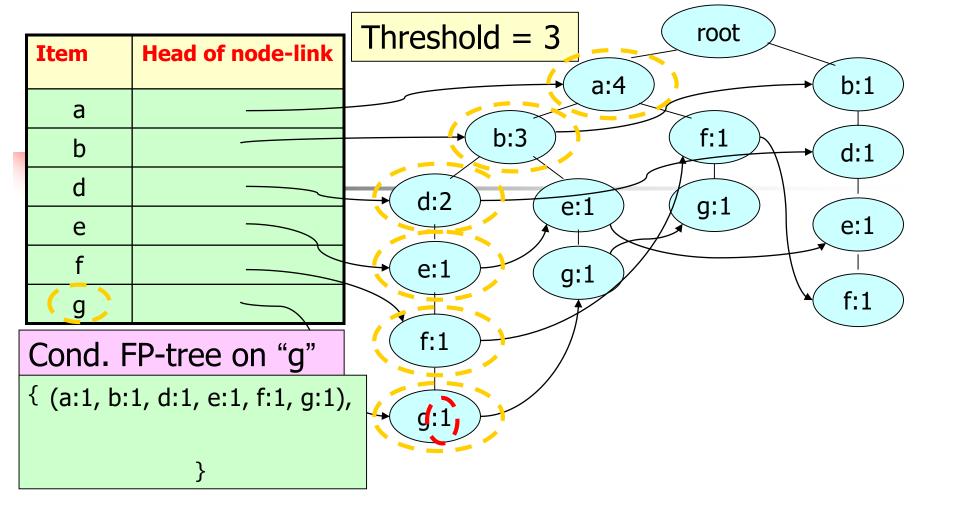
COMP1942

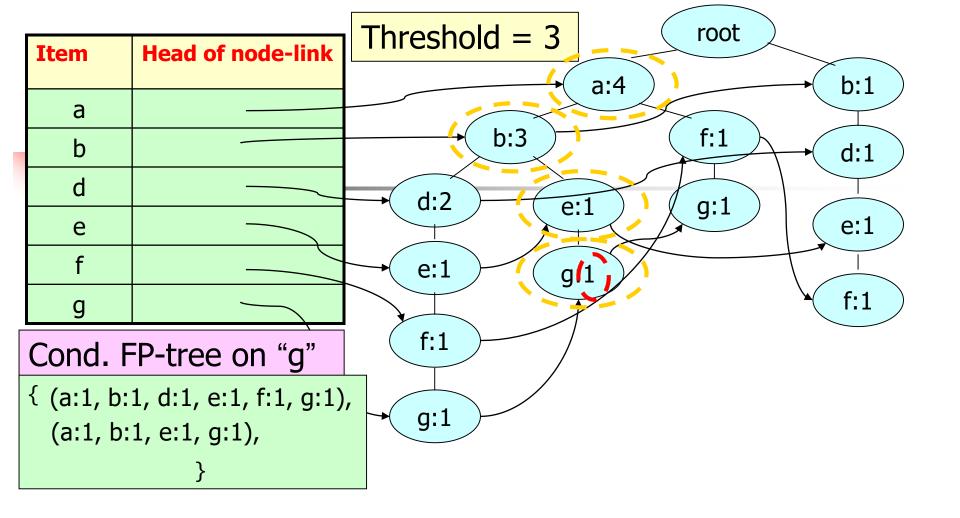
| TID | Items Bought | (Ordered) Frequent Items |
|-----|------------------------|--------------------------|
| 100 | a, b, c, d, e, f, g, h | a, b, d, e, f, g |
| 200 | a, f, g | a, f, g |
| 300 | b, d, e, f, j | b, d, e, f |
| 400 | a, b, d, i, k | a, b, d |
| 500 | a, b, e, g | a, b, e, g |

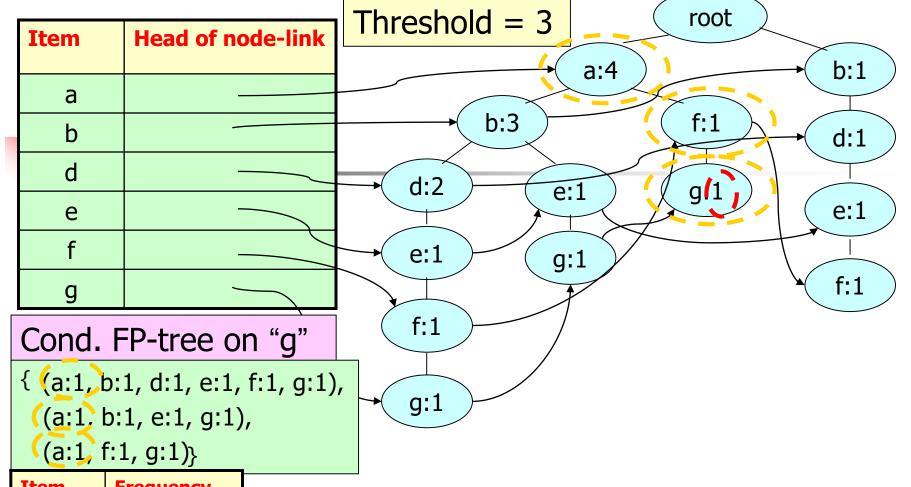




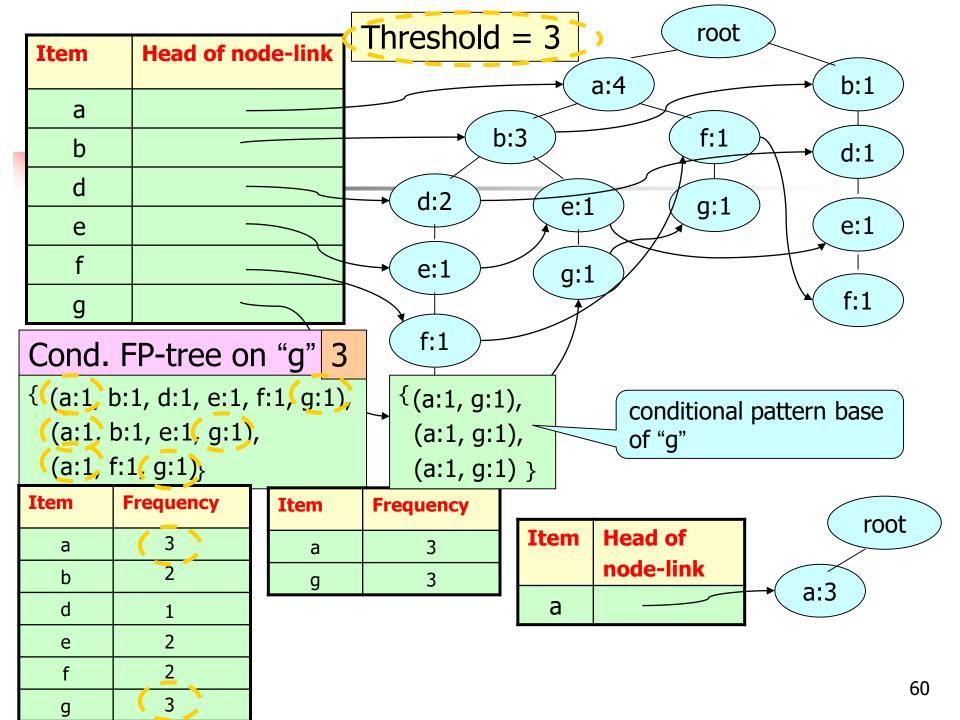


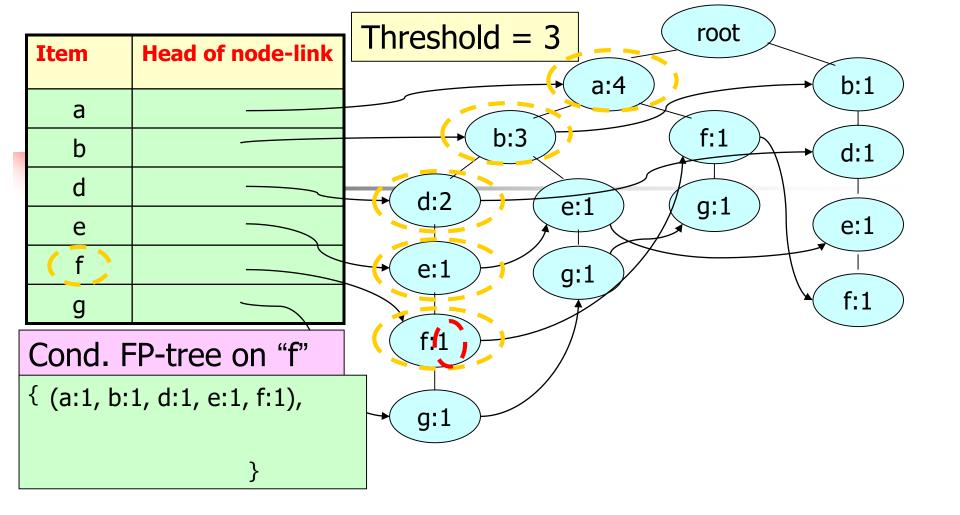


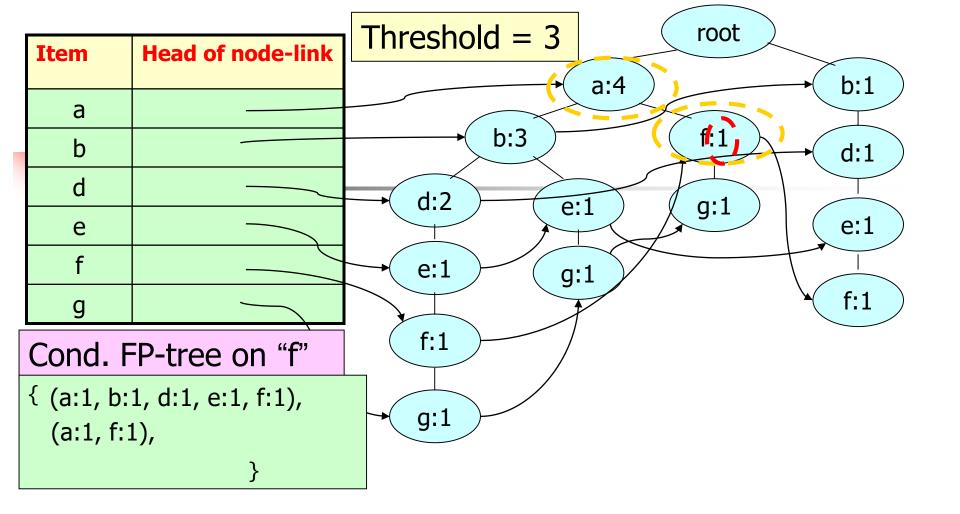


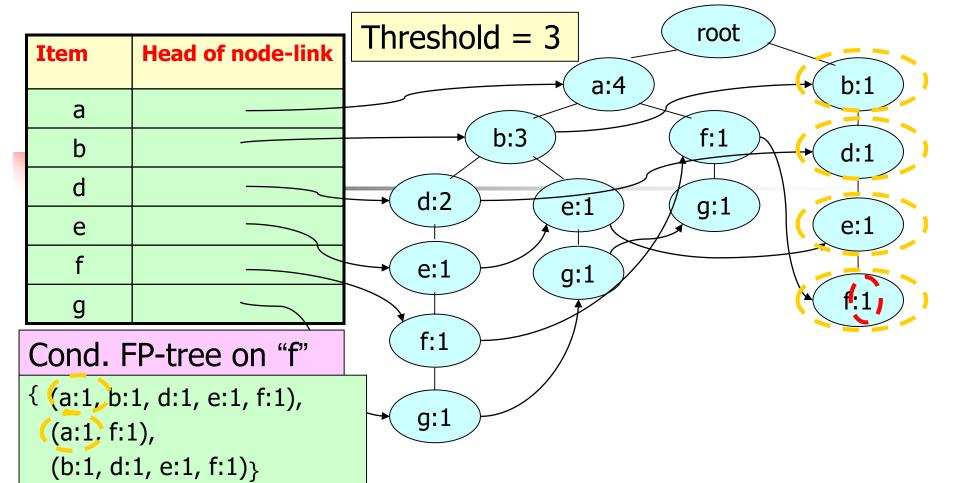


| Item | Frequency |
|------|-----------|
| а | 3 |
| b | 2 |
| d | 1 |
| е | 2 |
| f | 2 |
| g | 3 |

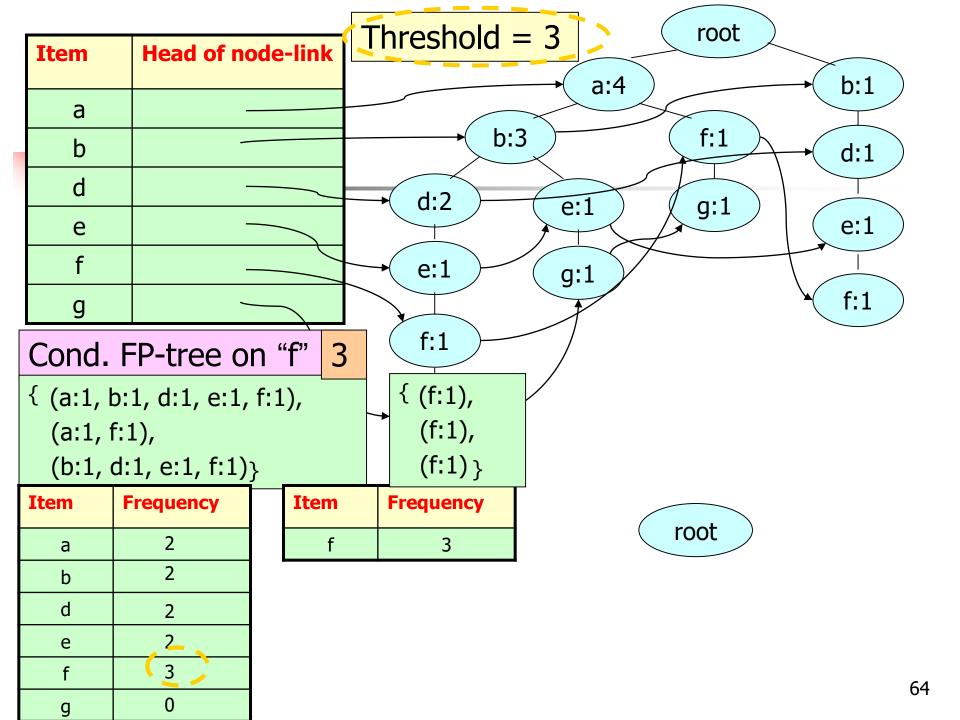


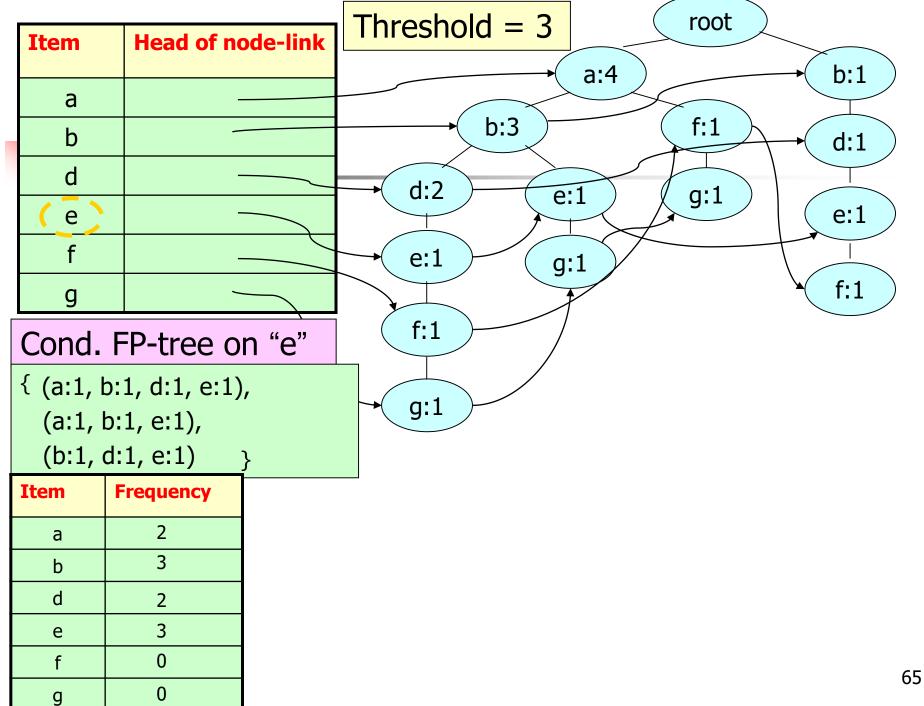


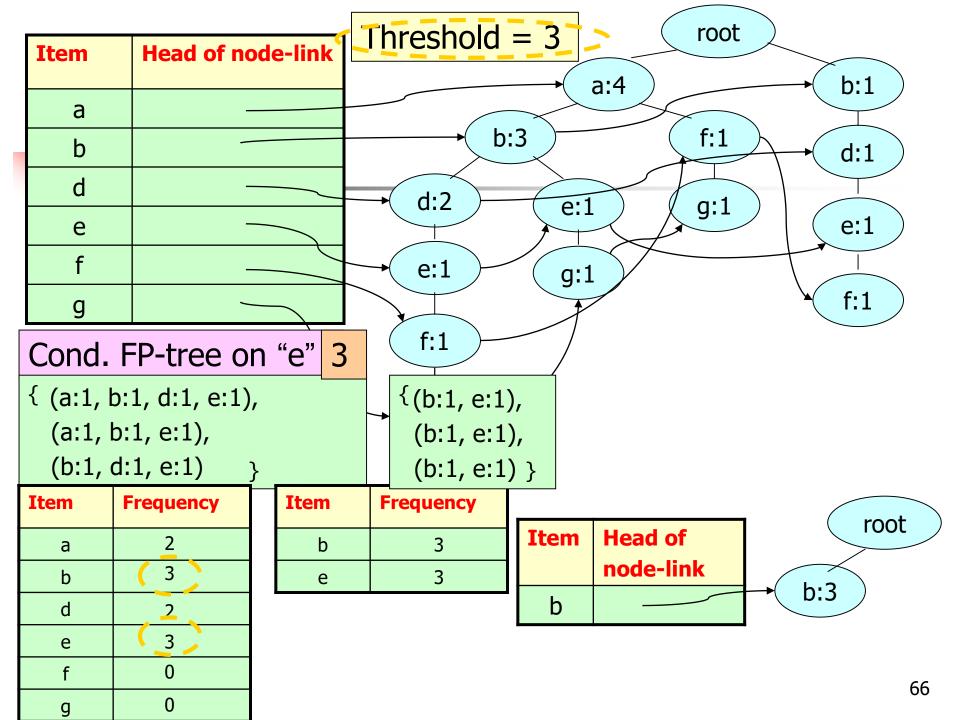


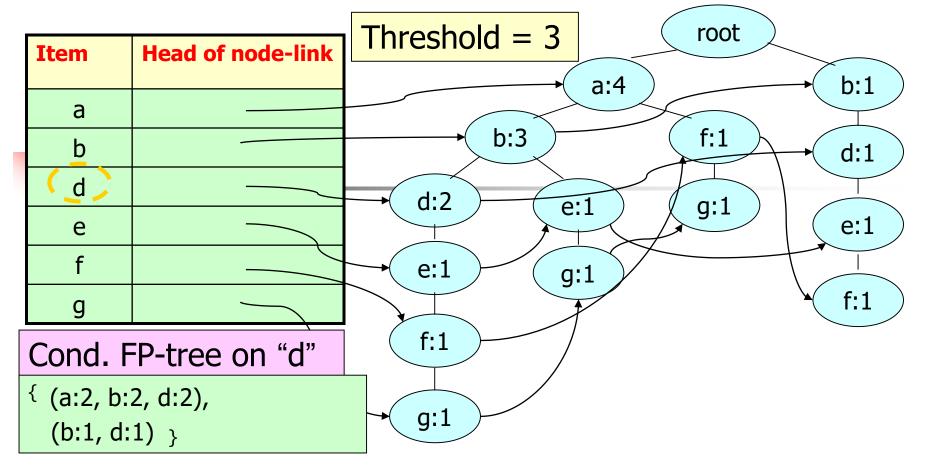


| Item | Frequency |
|------|-----------|
| а | 2 |
| b | 2 |
| d | 2 |
| е | 2 |
| f | 3 |
| g | 0 |

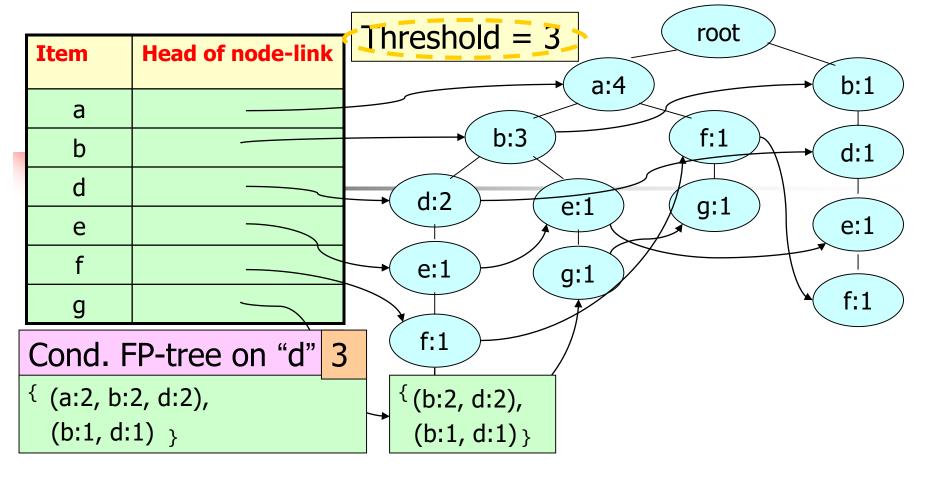








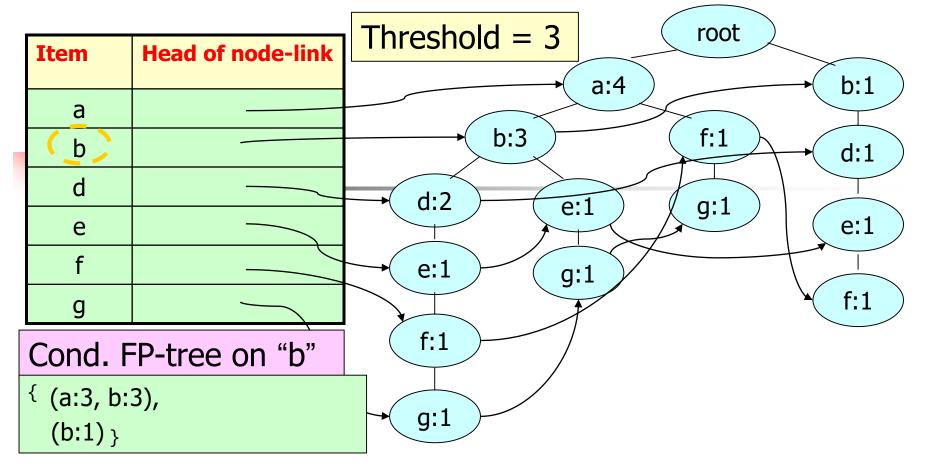
| Item | Frequency |
|------|-----------|
| а | 2 |
| b | 3 |
| d | 3 |
| е | 0 |
| f | 0 |
| g | 0 |



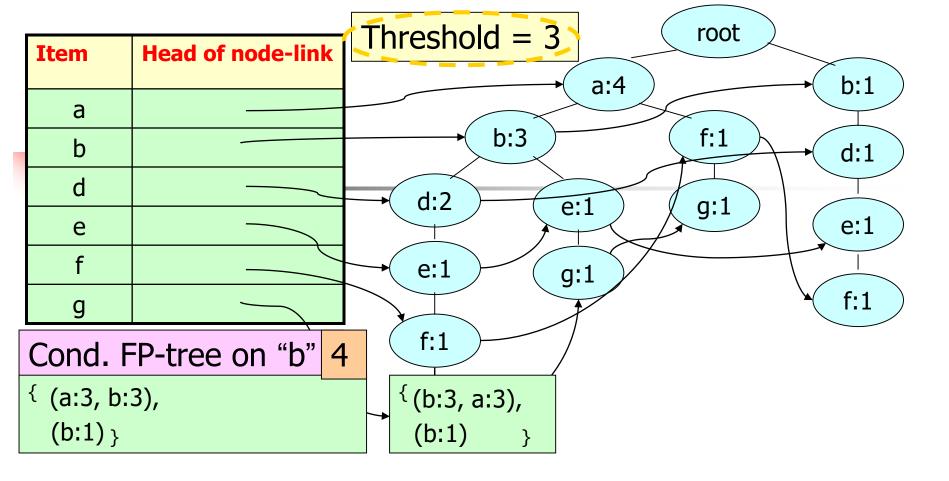
| Item | Frequency |
|------|-----------|
| а | 2 |
| b | 3 |
| d | 3 |
| е | 0 |
| f | 0 |
| g | 0 |

| Frequency |
|-----------|
| 3 |
| 3 |
| |

| Item | Head of node-link | root |
|------|-------------------|---------|
| b | | → (b:3) |



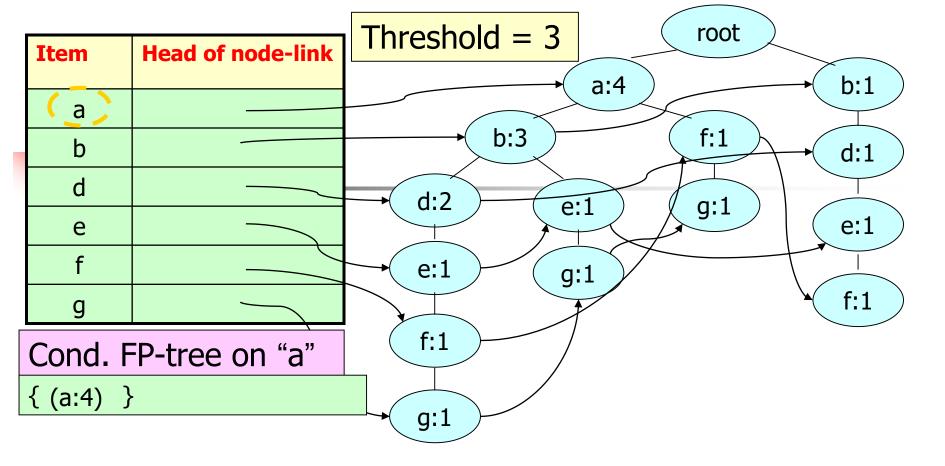
| Item | Frequency |
|------|-----------|
| а | 3 |
| b | 4 |
| d | 0 |
| е | 0 |
| f | 0 |
| g | 0 |



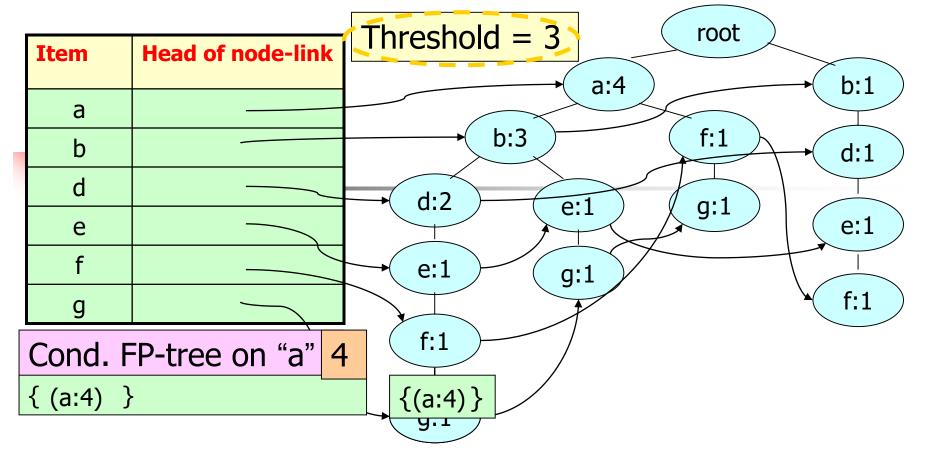
| Item | Frequency |
|------|-----------|
| а | 3 |
| b | 4 |
| d | 0 |
| е | 0 |
| f | 0 |
| g | 0 |

| Item | Frequency |
|------|-----------|
| b | 4 |
| a | 3 |

| Item | Head of node-link | root |
|------|-------------------|---------|
| а | | → (a:3) |



| Item | Frequency |
|------|-----------|
| а | 4 |
| b | 0 |
| d | 0 |
| е | 0 |
| f | 0 |
| g | 0 |



| Item | Frequency |
|------|-----------|
| а | 4 |
| b | 0 |
| d | 0 |
| e | 0 |
| f | 0 |
| g | 0 |

| Item | Frequency |
|------|-----------|
| а | 4 |

root

FP-tree

Step 1: Deduce the ordered frequent items. For items with the same frequency, the order is given by the alphabetical order.

Step 2: Construct the FP-tree from the above data

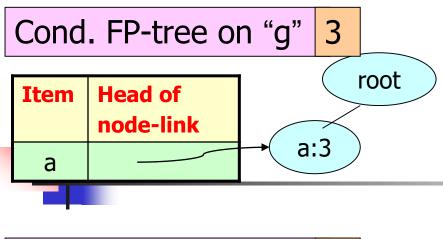
Step 3: From the FP-tree above, construct the FP-

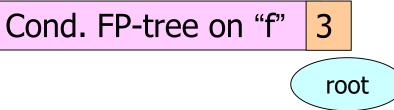
conditional tree for each item (or itemset).

Step 4: Determine the frequent patterns.

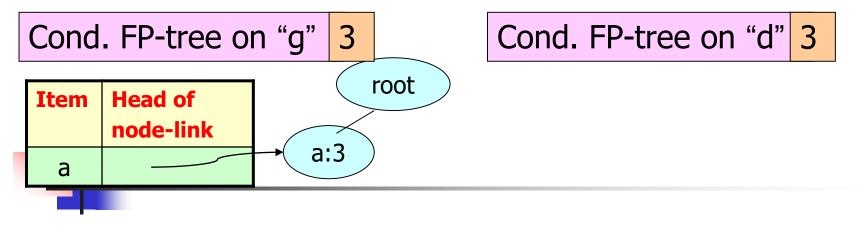
Cond. FP-tree on "g" 3

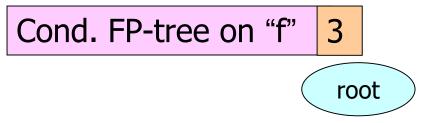


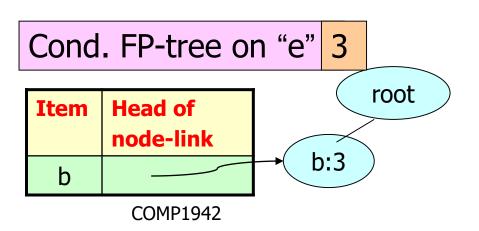


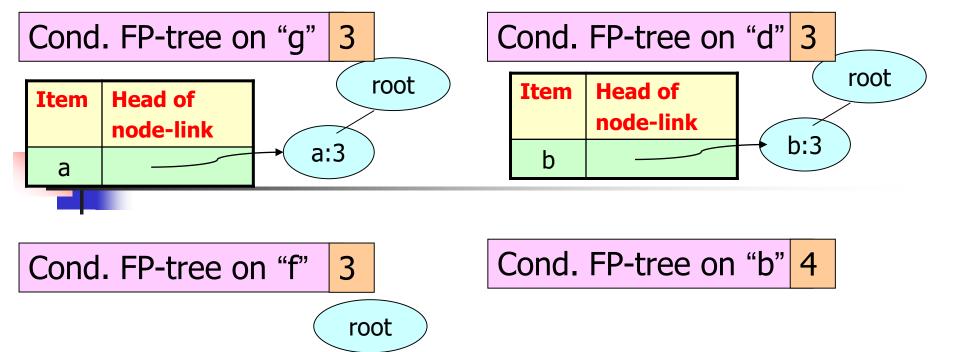


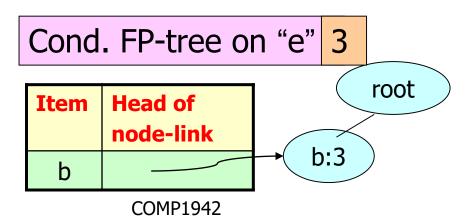
Cond. FP-tree on "e" 3

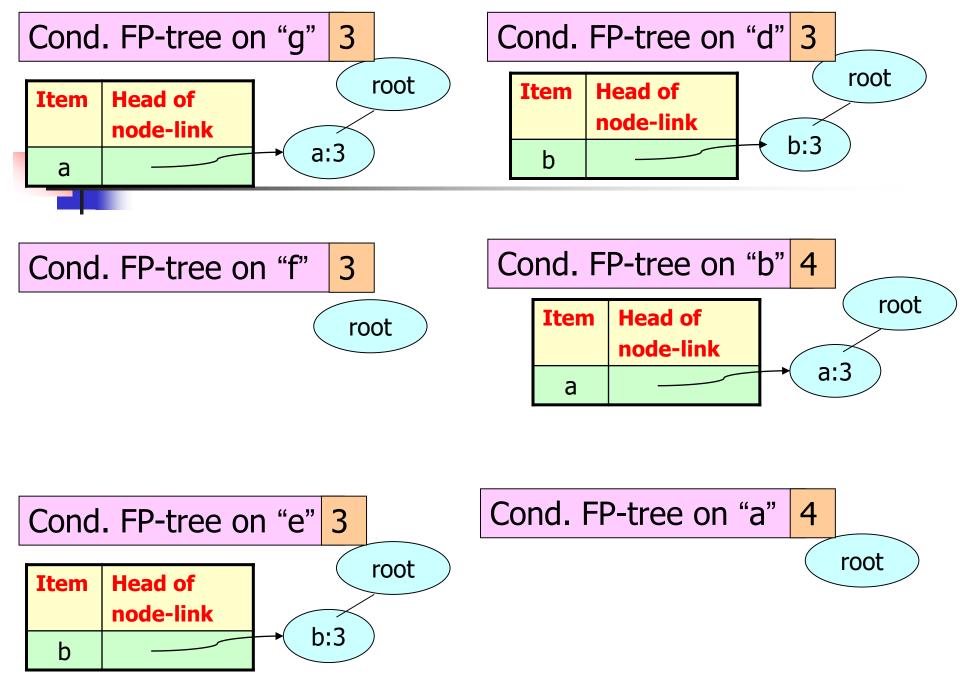






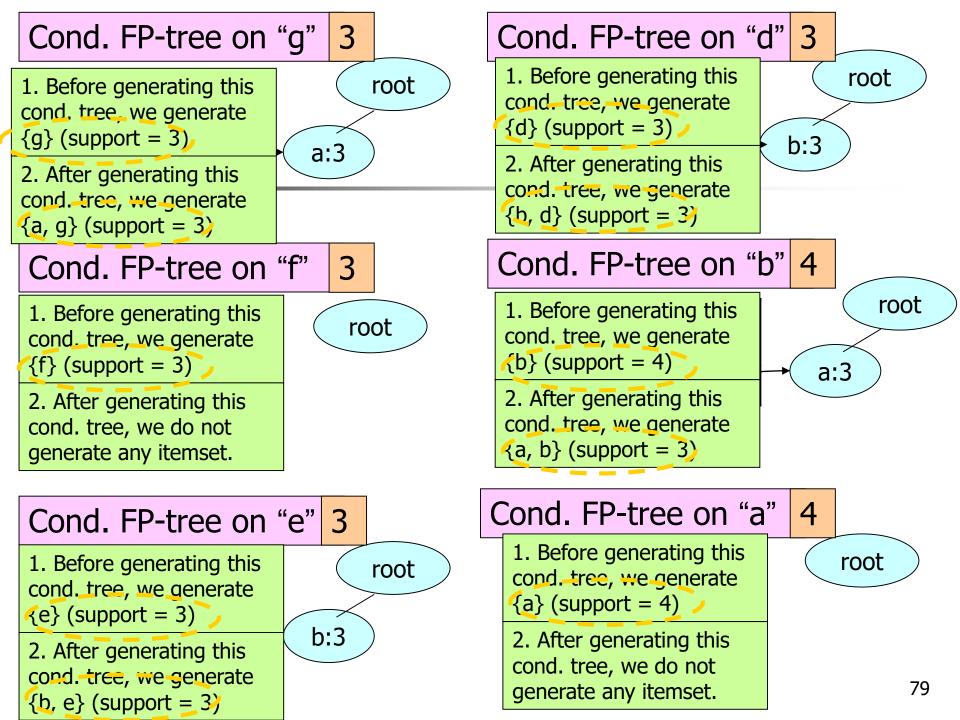






COMP1942

78



Complexity

- Complexity in building FP-tree
 - Two scans of the transactions DB
 - Collect frequent items
 - Construct the FP-tree
- Cost to insert one transaction
 - Number of frequent items in this transaction



Size of the FP-tree

The size of the FP-tree is bounded by the overall occurrences of the frequent items in the database



Height of the Tree

 The height of the tree is bounded by the maximum number of frequent items in any transaction in the database

Compression

- With respect to the total number of items stored,
 - is FP-tree more compressed compared with the original databases?

Details of the Algorithm

- Procedure FP-growth (Tree, α)
 - if Tree contains a single path P
 - for each combination (denoted by β) of the nodes in the path P do
 - generate pattern β U α with support = minimum support of nodes in β
 - else
 - for each a_i in the header table of Tree do
 - generate pattern $\beta = a_i U \alpha$ with support = a_i .support
 - construct β 's conditional pattern base and then β 's conditional FP-tree Tree $_{\beta}$
 - if Tree_{β} $\neq \emptyset$
 - Call FP-growth(Tree_β, β)

Outline

- Association Rule Mining
 - Problem Definition
- Algorithm
- How to use the data mining tool

How to use the data mining tool

Please see the next set of slides.

COMP1942 86