

Faculty of Computer Science and Engineering
Ho Chi Minh City University of Technology

Chapter 6

Association Rule

TRAN MINH QUANG

quangtran@hcmut.edu.vn

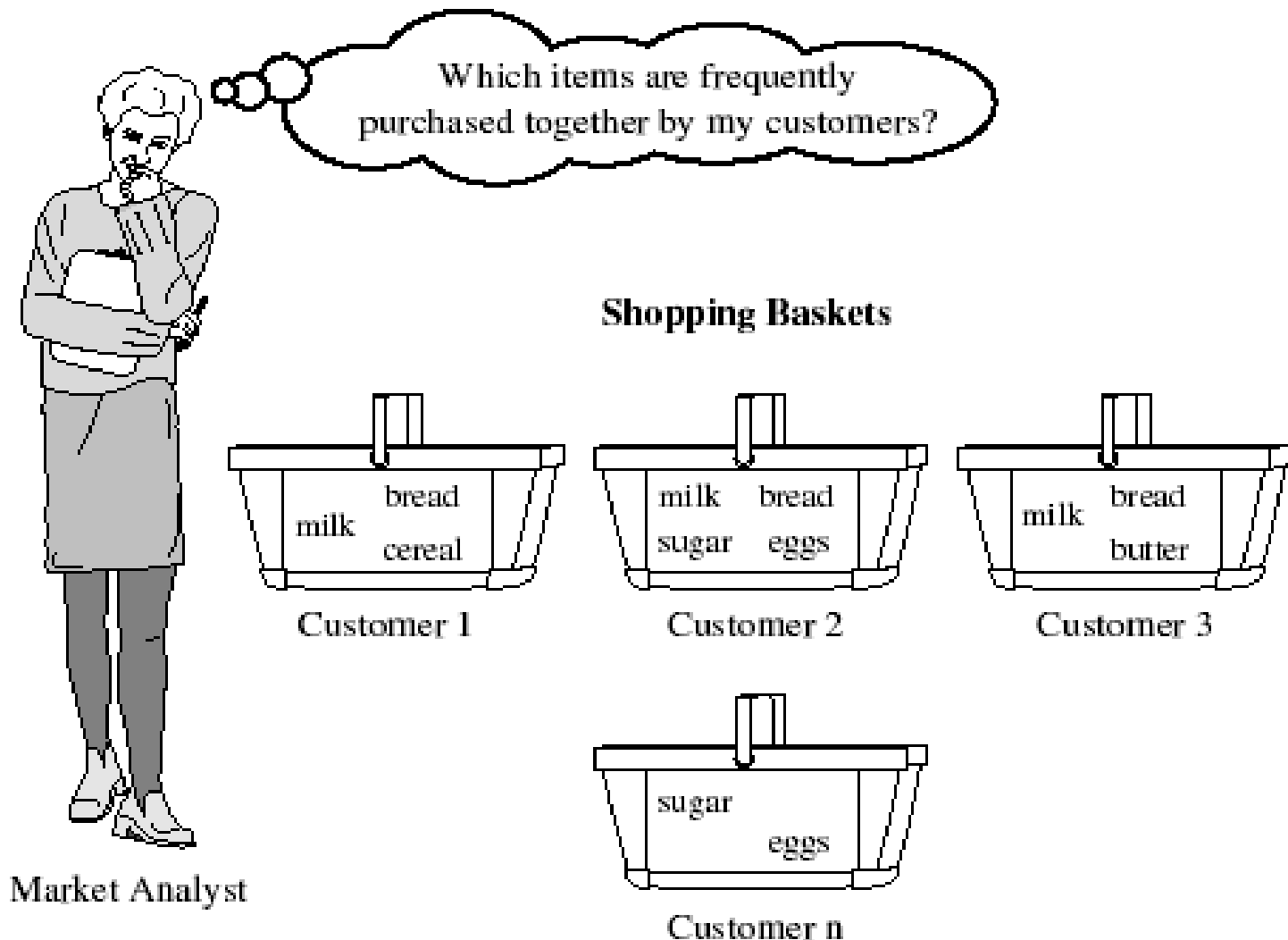
<http://researchmap.jp/quang>

1

CONTENT

1. Situations
2. Overview on association rules
3. Association rule representation
4. Frequent itemset mining
 1. Apriori
 2. FP-Growth
5. Mining association rules from frequent itemsets
6. Mining association rules based on constraints
7. Correlation analysis
8. Summary

1. SITUATION 1 – BASKET ANALYSIS




1. SITUATION 2 – RECOMMENDATION

Amazon.com: Information Systems: Foundation of E-Business (4th Edition) (9780130617736): Steven - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Search Favorites Media

Customers Who Bought This Item Also Bought

Image	Product Title	Author	Rating	Price
	Management Information Systems: Managing the...	by Kenneth C. Laudon	★★★★☆ (5)	
	The Internet Book: Everything You Need to Know...	by Douglas E. Comer	★★★★☆ (10)	\$47.53
	Project Management with MS Project CD + Student...	by Erik W. Larson	★★★★☆ (20)	\$141.85

Editorial Reviews

Product Description

Emphasizes the essential role of information systems in the works systems through which today's businesses operate. For professionals in the field of information systems.

From the Back Cover

Internet

1. SITUATION 2 – RECOMMENDATION

Đặc Nhân Tâm - Cuốn Sách Hay Nhất Của Mọi Thời Đại Đưa Bạn Đến Thành Công - Sách Vinabook.com - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites Media Print Mail News RSS Feeds

NHÂN TÂM
How to Win Friends & Influence People
Cuốn sách hay nhất của mọi thời đại đưa bạn đến thành công
NEW EDITION

★★★★★ (11 người đánh giá)

Số trang: 320
Kích thước: 14.5x20.5 cm
Trọng lượng: 370 gram
(Chi tiết về phí vận chuyển)

Hình thức bìa: Bìa mềm
Ngày xuất bản: 09 - 2008
Số lần xem: 87539

[In trang này](#)
[Gửi cho bạn bè](#)

Giá bìa: 48.000 VNĐ
Giá bán: **48.000 VNĐ**
Giảm giá: (0%)

Vietnam đồng

Xếp hạng: 12 (trong những cuốn Sách bán chạy)

Sách nên mua kèm với sách này
Nên mua cuốn sách trên cùng với cuốn **Quảng Gánh Lo Đi Và Vui Sống - Những Ý Tưởng Tuyệt Vời Để Sống Thanh Thản Và Hạnh Phúc** - Nhà Trẻ

ĐẶC NHÂN TÂM + **Quảng gánh lo đi & Vui sống**

Giá bán tất cả: 96.000 VNĐ

[Thêm tất cả vào giỏ hàng](#)

Khách hàng mua cuốn sách trên cũng từng mua một trong những cuốn sách dưới đây

Trang: 1 / 10

Người Giỏi Không Phải Là Người Làm ... Tác giả: Donna M. Genett
Giá bán: 24.000 VNĐ

Quảng gánh lo đi & Vui sống Tác giả: Dale Carnegie
Giá bán: 48.000 VNĐ

Lập Bản Đồ Tư Duy - Công Cụ Tư Duy ... Tác giả: Tony Buzan
Giá bán: 24.000 VNĐ

Nguyên Lý 80/20 - Bí Quyết Làm Ít Được Nhiều ... Tác giả: Richard Koch
Giá bán: 70.000 VNĐ

Bài Giảng Cuối Cùng - The Last Lecture ... Tác giả: Randy Pausch
Giá bán: 58.000 VNĐ

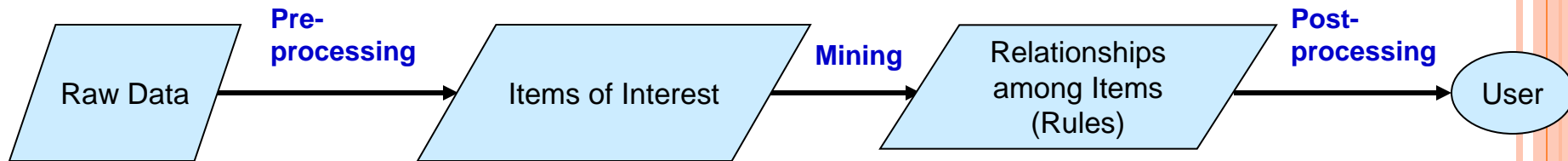
Done Internet

1. SITUATION ...

- Basket data analysis
- Cross-marketing, recommendation,...
- Catalog design
- Classification and clustering with frequent patterns
- ...

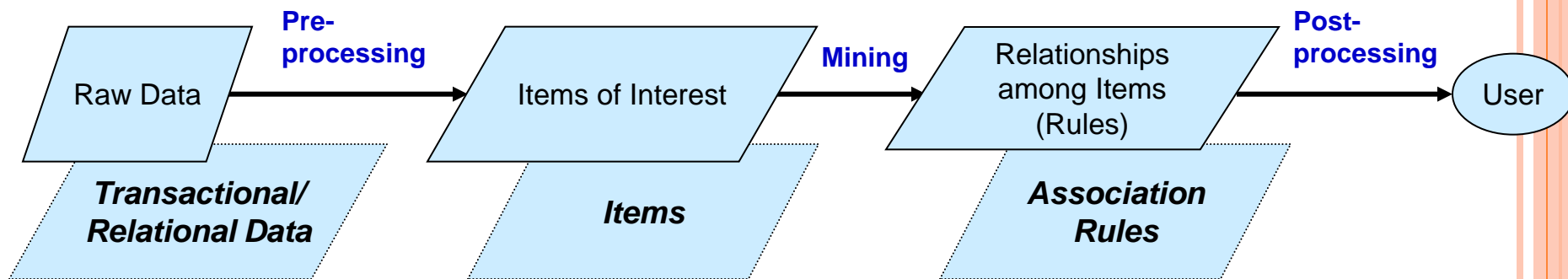
2. OVERVIEW ON ASSOCIATION RULE MINING

- Association rule mining process



2. OVERVIEW ON ASSOCIATION RULE MINING

○ Association rule mining process



Transaction	Items_bought
2000	A, B, C
1000	A, C
4000	A, D
5000	B, E, F
...	

A, B, C, D, F,
...

$A \rightarrow C$ (50%, 66.6%)
...

Basket-based analysis

2. OVERVIEW ON ASSOCIATION RULE MINING

- Basic concepts
 - Item
 - Itemset
 - Transaction
 - Association and association rules
 - Support
 - Confidence
 - Frequent itemset
 - Strong association rules

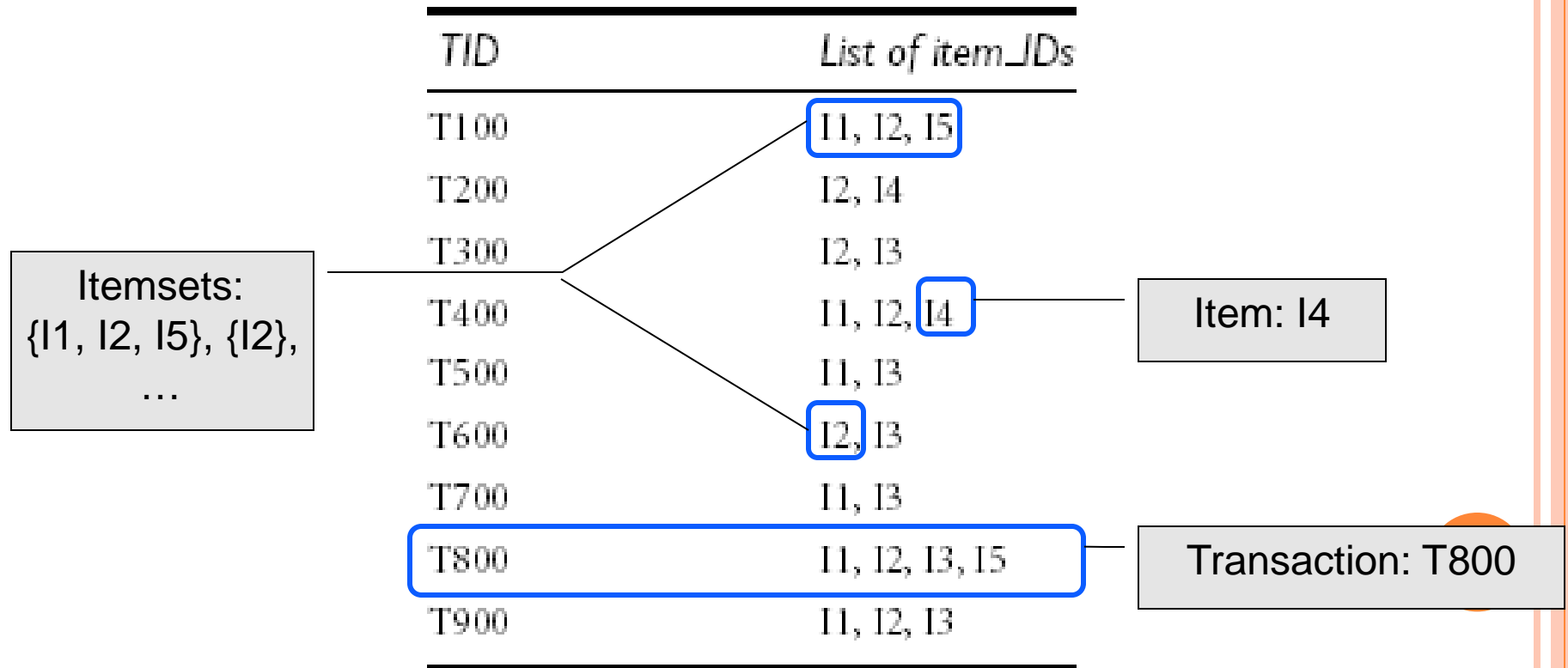
2. OVERVIEW ON ASSOCIATION RULE MINING

- AllElectronics dataset (after preprocessing)

<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

2. OVERVIEW ON ASSOCIATION RULE MINING

- AllElectronics dataset (after preprocessing)



2. OVERVIEW ON ASSOCIATION RULE MINING

○ Basic concepts

- Item: pattern, sample, object of interest
- $J = \{I_1, I_2, \dots, I_m\}$: a set of all m items in the dataset
- Itemset
 - A set of items
 - k-itemset: itemset of k items
- Transaction
 - A record in a transactional dataset
 - A set of T items in the same transaction/record

2. OVERVIEW ON ASSOCIATION RULE MINING

○ Basic concepts

- Association vs. association rule
 - Association: occurrence of items in the same transaction(s)
 - Represents the linkage/association between items/itemsets
- Association rule: a criteria/rule for the association between itemsets
 - Represents the linkage (with some conditions) between itemsets
 - Let A and B be itemsets, an association rule between A and B is $A \rightarrow B$ (B occurs in the condition that A occurs)

2. OVERVIEW ON ASSOCIATION RULE MINING

○ Basic concept

- Support: measure the occurrence frequency of items/itemsets
- Minimum support threshold (MinSup): is the minimum support value defined by user to validate the “support” of a frequent itemset/association
- Confidence: to measure the occurrence frequency of an itemset upon the occurrence of another itemset
- Minimum confidence threshold (MinConf): is the minimum confidence value defined by user to validate the “confidence” of an association rule.

2. OVERVIEW ON ASSOCIATION RULE MINING

○ Basic concepts

- Frequent itemset: is an itemset whose support satisfies the minimum support threshold

A is a frequent itemset iff

$$\text{support}(A) \geq \text{MinSup}$$

- Strong association rule: is an association rule whose support and confidence satisfy MinSup and MinConf
 - Give $A \rightarrow B$, where A and B are itemsets
 - $A \rightarrow B$ is a strong association rule iff

$$\begin{aligned} \text{support}(A \rightarrow B) &\geq \text{MinSup} \text{ and} \\ \text{confidence}(A \rightarrow B) &\geq \text{MinConf} \end{aligned}$$

2. OVERVIEW ON ASSOCIATION RULE MINING

- Association rule categories
 - Boolean association rule vs. quantitative association rule
 - Single-dimensional association rule vs. multidimensional association rule
 - Single-level association rule vs. multilevel association rule
 - Association rule vs. correlation rule

2. OVERVIEW ON ASSOCIATION RULE MINING

- Boolean vs. quantitative association rules
 - Boolean association rule: represents the association of occurrence/absence of itemsets

Computer \rightarrow *Financial_management_software*
[support=50%, confidence=60%]

- Quantitative association rule: represents the association of between quantitative items/features

Age(X, “30..39”) \wedge Income(X, “42K..48K”) \rightarrow buys(X, high resolution TV)

[support=50%, confidence=60%]

2. OVERVIEW ON ASSOCIATION RULE MINING

- Single-dimensional vs. multidimensional association rules
 - Single-dimensional association rule: is the rule that relates to only one dimension of items

Buys(X, “computer”) \rightarrow *Buys*(X,
“financial_management_software”)

- Multidimensional association rule: is the rule that relates to multiple dimensions of items

Age(X, “30..39”) \rightarrow *Buys*(X, “computer”)

2. OVERVIEW ON ASSOCIATION RULE MINING

- Single-level vs. multilevel association rules
 - Single-level association rule: is the rule that relates to items/features in one level of abstraction

$$\text{Age}(X, \text{"30..39"}) \rightarrow \text{Buys}(X, \text{"computer"})$$
$$\text{Age}(X, \text{"18..29"}) \rightarrow \text{Buys}(X, \text{"camera"})$$

- Multilevel association rule: is the rule that relates to items/features in multiples levels of abstraction

$$\text{Age}(X, \text{"30..39"}) \rightarrow \text{Buys}(X, \text{"laptop computer"})$$
$$\text{Age}(X, \text{"30..39"}) \rightarrow \text{Buys}(X, \text{"computer"})$$

2. OVERVIEW ON ASSOCIATION RULE MINING

- Association rule vs. correlation rule
 - Association rule: strong association rule $A \rightarrow B$ is that satisfies minimum support threshold and minimum confidence threshold
 - Correlation rule: strong correlation rules $A \rightarrow B$ is that satisfies correlation conditions in statistics (using different statistical correlation measures)

3. ASSOCIATION RULE REPRESENTATION
























Rule form:

$$A \rightarrow B [support, confidence]$$

Where,

- **A, B are frequent itemsets**
- $Support(A \rightarrow B) = Support(A \cup B) \geq min_sup$
- $Confidence(A \rightarrow B) = P(B | A) = Support(A \cup B) / Support(A) \geq min_conf$

3. ASSOCIATION RULE REPRESENTATION



Support



Sup = 1



Sup = 1

Confidence

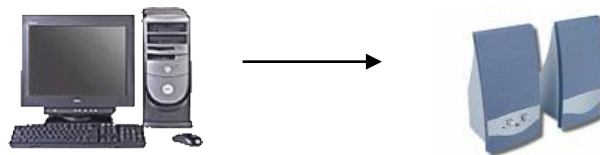
$$C(\text{Computer} \rightarrow \text{Speakers}) = \frac{4}{5}$$

4. MINING FREQUENT ITEMSETS/PATTERNS

- 2 steps in association rule mining
 - Finding frequent itemsets/patterns



- Mining association rules from frequent itemsets/patterns



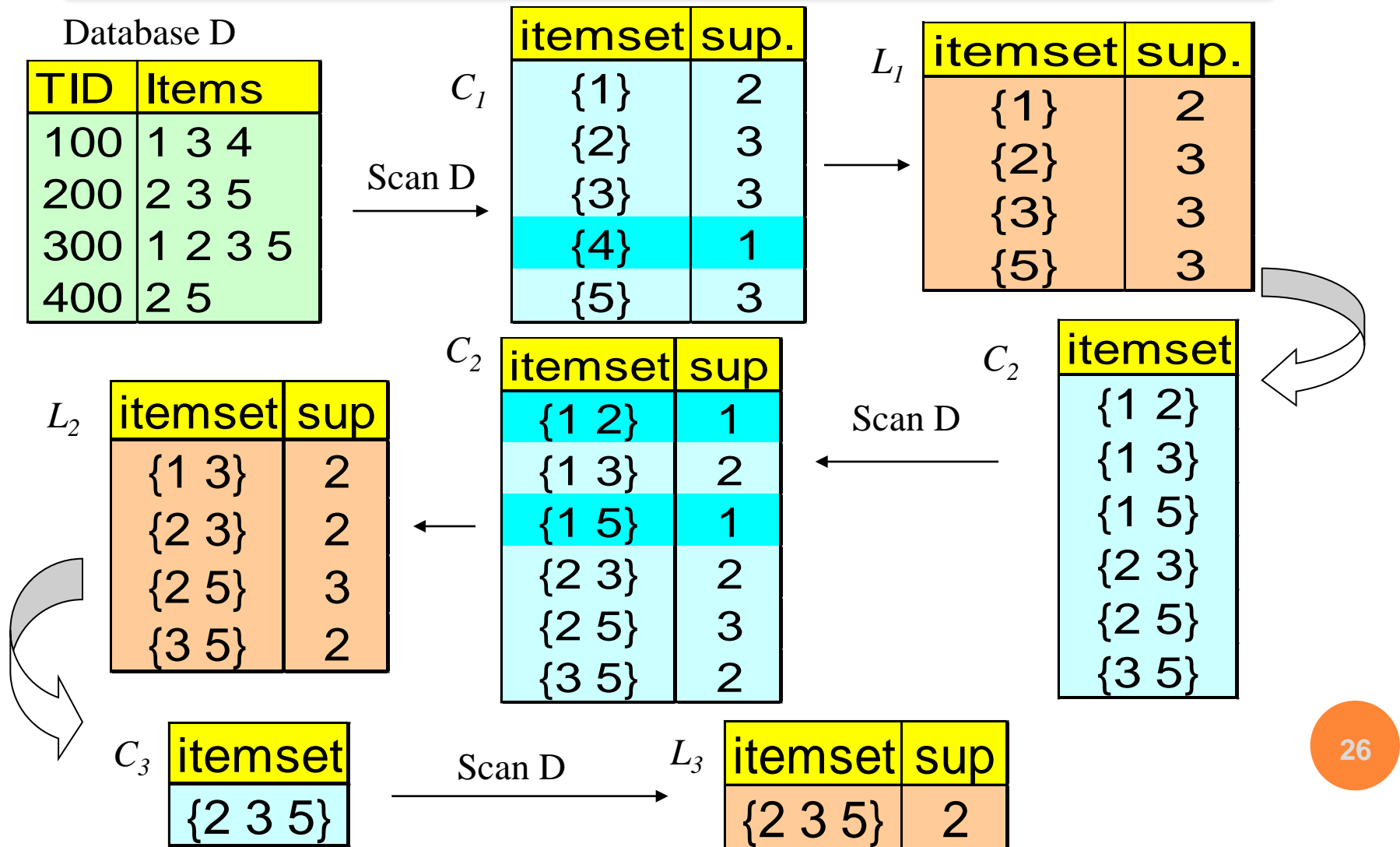
4. MINING FREQUENT ITEMSETS/PATTERNS

- Apriori method: mining frequent itemsets with candidate itemsets using prior knowledge
 - R. Agrawal, R. Srikant. Fast algorithms for mining association rules. In VLDB 1994, pp. 487-499
- FP-Growth: using FP-tree
 - J. Han, J. Pei, Y. Yin. Mining frequent patterns without candidate generation. In SIGMOD 2000, pp. 1-12

4.1. MINING FREQUENT ITEMSETS - APRIORI

- Use the prior knowledge about the characteristics of frequent itemsets
- Iterate the searching process to find frequent itemsets at each level (level-wise search)
 - $k+1$ -itemsets: generated from k -itemsets
 - At each level, check all the data set to identify frequent itemsets
- Apriori property (to reduce the search space): All subsets of a frequent itemset are frequent itemsets
- Anti-monotone: If X is NOT a frequent itemset then **neither** does $\{X \cup Y\}$

4.1. MINING FREQUENT ITEMSETS - APRIORI



4.1. MINING FREQUENT ITEMSETS - APRIORI

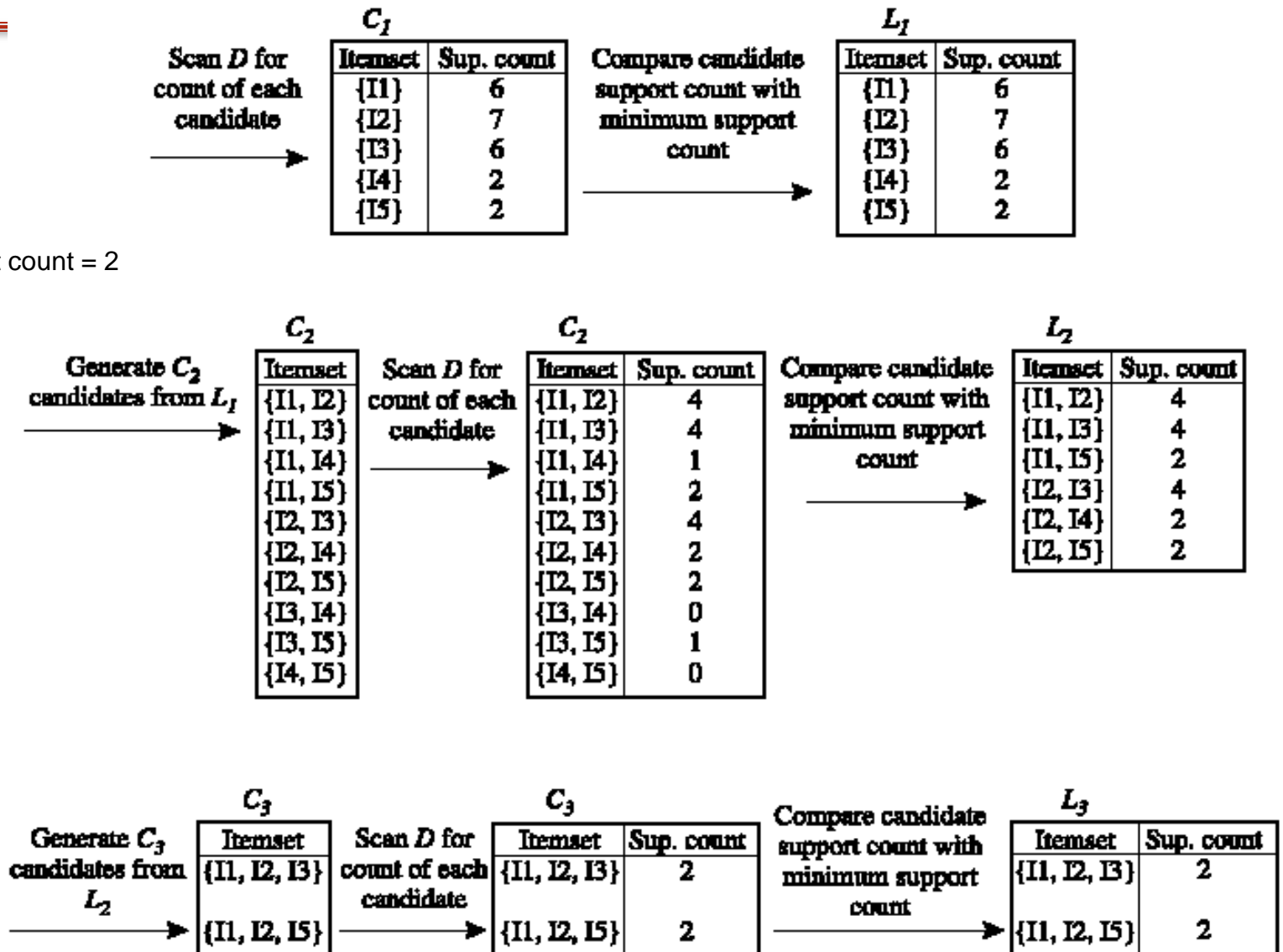
- AllElectronics dataset (after preprocessing)

<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

4.1. MINING FREQUENT ITEMSETS - APRIORI

min_sup = 2/9

minimum support count = 2



4.1. MINING FREQUENT ITEMSETS - APRIORI

- Pseudo-code:

C_k : Candidate itemset of size k

L_k : frequent itemset of size k

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

While ($L_k \neq \emptyset; k++$) **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 $\cup_k L_k$;

4.1. MINING FREQUENT ITEMSETS - APRIORI

- How to generate candidates?
 - Step 1: self-joining L_k
 - Step 2: pruning (**Antimonotone**)
- How to count supports of candidates?
- 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
 - $C_4 = \{abcd\}$

4.1. MINING FREQUENT ITEMSETS - APRIORI

- Suppose the items in L_{k-1} are listed in an order
- Step 1: self-joining L_{k-1}
 - insert into C_k
 - select $p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$
 - from $L_{k-1} p, L_{k-1} q$
 - where $p.item_1=q.item_1, \dots, p.item_{k-2}=q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$
- Step 2: pruning
 - forall *itemsets* c in C_k do
 - forall *(k-1)-subsets* s of c do
 - if (s is not in L_{k-1}) then delete c from C_k

4.1. MINING FREQUENT ITEMSETS - APRIORI

- Apriori's main characteristics
 - Create many candidates for frequent itemsets
 - 10^4 frequent 1-itemsets \rightarrow more than 10^7 ($\approx 10^4(10^4 - 1)/2$) 2-itemsets as candidates
 - For each k-itemset, it examines $2^k - 1$ itemsets of candidates
 - Scan the original data (D) many times for calculating the supports
 - Huge cost when the size of candidate itemsets increase
 - If k-itemsets is mined D is scanned k+1 times

4.1. MINING FREQUENT ITEMSETS - APRIORI

○ Apriori's improvement techniques

- Techniques based on hashing: An k -itemset whose hashing bucket count is smaller than minimum support threshold is **not** a frequent itemset.
- Reduce the scan time: A transaction that does not contain **any** frequent k -itemset are not needed to be examined at later level (i.e., for identify $k+1$ -itemset).
- Partitioning: An itemset must be frequent in at least one partition to be considered as frequent in the whole dataset.
- Sampling: Mine the frequent itemsets from samples (via sampling) with a smaller support threshold. It is necessary to have a method for validating the completeness.
- Dynamic itemset counting: Add an itemset to the candidates only iff all of its subsets are **predicted** to be “frequent”.

4.2. MINING FREQUENT ITEMSETS - FPGROWTH

- Compact the data into a Frequent Pattern tree (FP-tree)
 - Reduce the size of the examined dataset -> reduce the mining cost
 - **Infrequent items are discarded early in the process**
- Divide-and-conquer approach
 - The mining process is divided into smaller tasks
 1. Build the FP-tree
 2. Mine frequent itemsets using the FP-tree
- Avoid creating large candidate itemsets
 - Each mining task examines a small portion of the dataset

4.2. MINING FREQUENT ITEMSETS - FPGROWTH

- 1. Building the FP-tree
 - 1.1. Scan D to find *frequent 1-itemsets*
 - 1.2. Order *frequent 1-itemsets* in the descending of support count (frequency)
 - 1.3. Scan D again (2nd scan), build the FP-tree
 - Initiate with root node, label it with “null” or {}
 - Each transaction/tuples will correspond to a branch in the FP-tree
 - Each node contains an item of the transaction
 - Items in a transaction are sorted in descending order
 - Each node associates with support count of the corresponding item
 - Transactions contain the same items will create shared prefix branches

4.2. MINING FREQUENT ITEMSETS - FPGROWTH

Input:

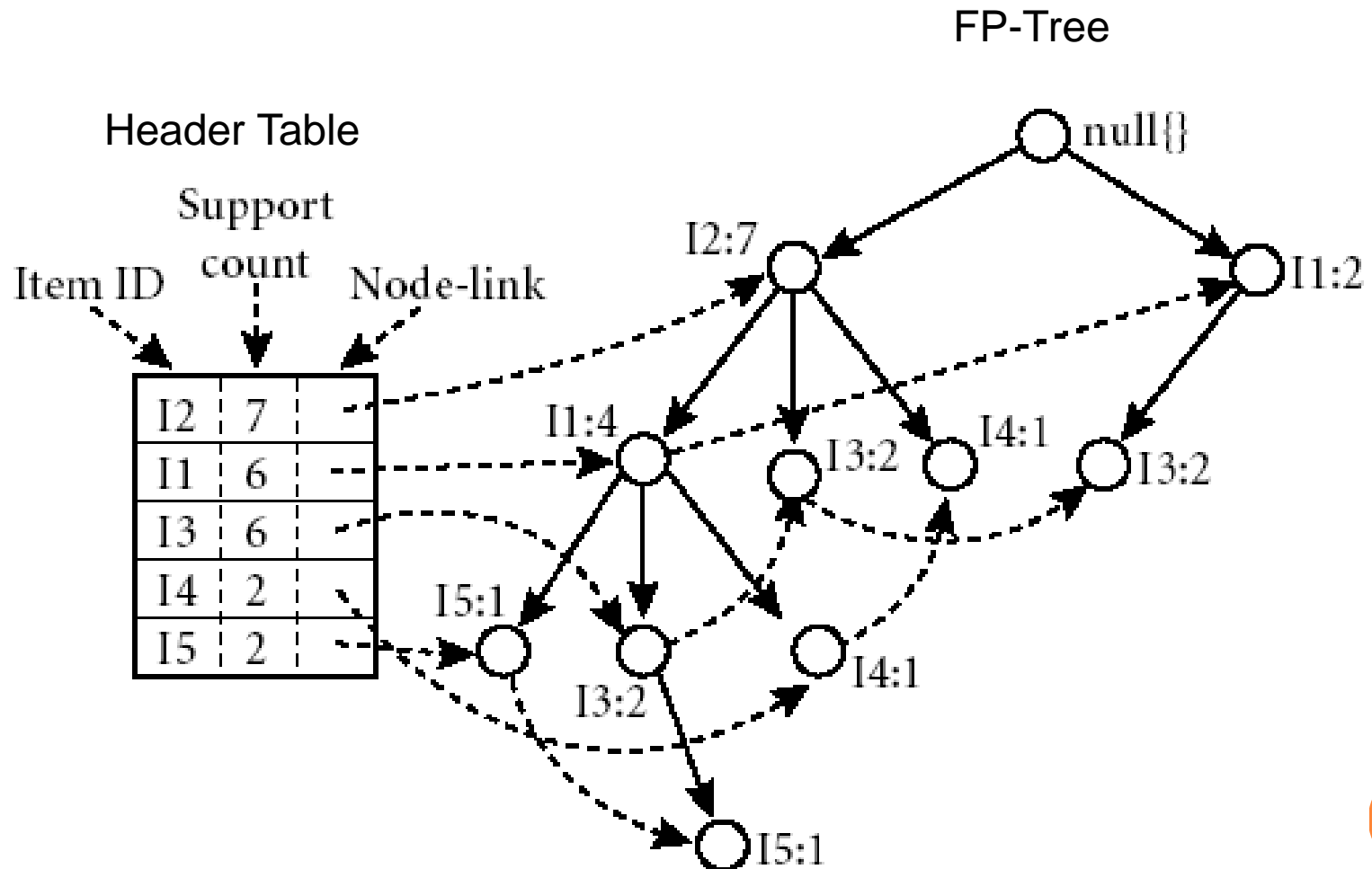
- D , a transaction database;
- min_sup , the minimum support count threshold.

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.” For each transaction $Trans$ in D do the following. Select and sort the frequent items in $Trans$ according to the order of L . Let the sorted frequent item list in $Trans$ be $[p|P]$, where p is the first element and P is the remaining list. Call $insert_tree([p|P], T)$, which is performed as follows. 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 $FP_growth(FP_tree, null)$, which is implemented as follows.

4.2. MINING FREQUENT ITEMSETS - FPGROWTH



4.2. MINING FREQUENT ITEMSETS - FPGROWTH

- 2. Mining frequent itemsets from the FP-tree
 - 2.1. Generate conditional pattern base for each node of the FP-tree
 - Accumulate the prefix paths with frequency of that node
 - 2.2. Create conditional FP-tree from each conditional pattern base
 - Accumulate frequency for each item in each base
 - Build conditional FP-tree for frequent items of that base
 - 2.3. Mine the conditional FP-tree for frequent itemsets recursively
 - If conditional FP-tree has a single path then list all itemsets

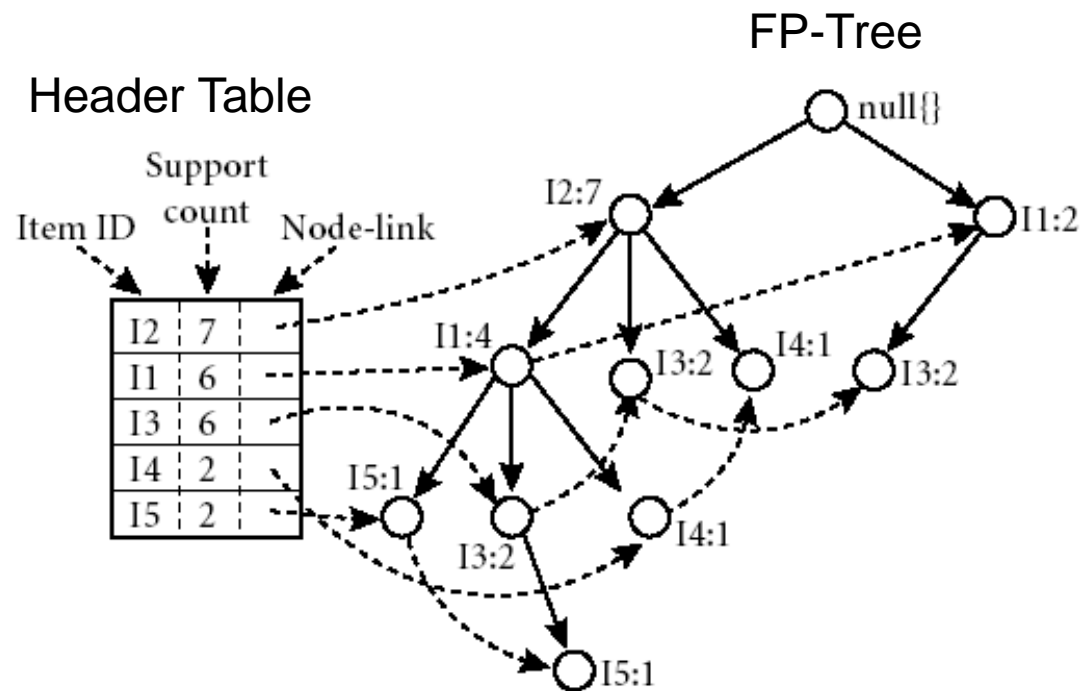
4.2. MINING FREQUENT ITEMSETS - FPGROWTH

- FP-Growth algorithm

procedure $\text{FP_growth}(Tree, \alpha)$

- (1) if $Tree$ contains a single path P then
- (2) for each combination (denoted as β) of the nodes in the path P
- (3) generate pattern $\beta \cup \alpha$ with *support_count* = *minimum support count of nodes in β* ;
- (4) else for each a_i in the header of $Tree$ {
- (5) generate pattern $\beta = a_i \cup \alpha$ with *support_count* = $a_i.\text{support_count}$;
- (6) construct β 's conditional pattern base and then β 's conditional FP-tree $Tree_\beta$;
- (7) if $Tree_\beta \neq \emptyset$ then
- (8) call $\text{FP_growth}(Tree_\beta, \beta)$; }

4.2. MINING FREQUENT ITEMSETS - FPGROWTH



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	$\{\{I2: 4\}\}$	$\langle I2: 4 \rangle$	$\{I2, I1: 4\}$

4.2. MINING FREQUENT ITEMSETS – FPGROWTH (ANOTHER EXAMPLE)

<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, w}	{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}

min_support = 3

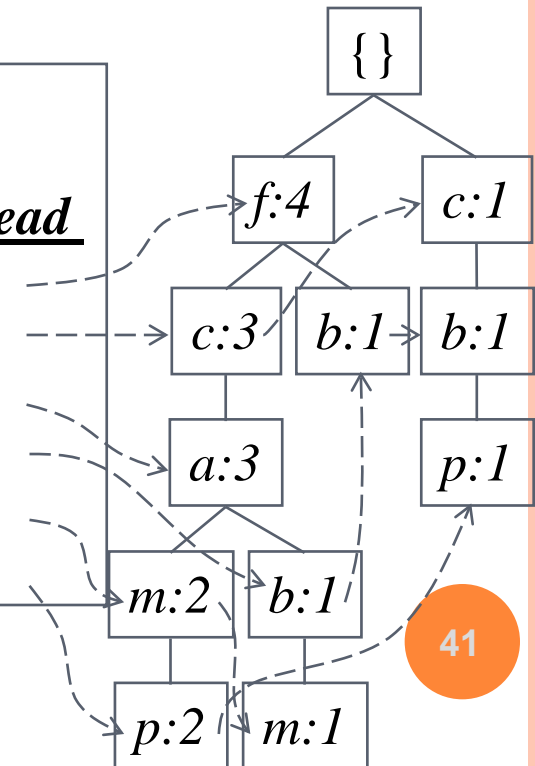
1. Scan DB once, find frequent 1-itemset (single item pattern)
2. Sort frequent items in frequency descending order, f-list
3. Scan DB again, construct FP-tree

Header Table

Item frequency head

<i>f</i>	4
<i>c</i>	4
<i>a</i>	3
<i>b</i>	3
<i>m</i>	3
<i>p</i>	3

F-list=f-c-a-b-m-p



4.2. FP-GROWTH EXAMPLE:

PARTITION PATTERNS AND DATABASES

- Frequent patterns can be partitioned into subsets according to f-list

f-list=f-c-a-b-m-p

Patterns containing p

Patterns having m but no p

...

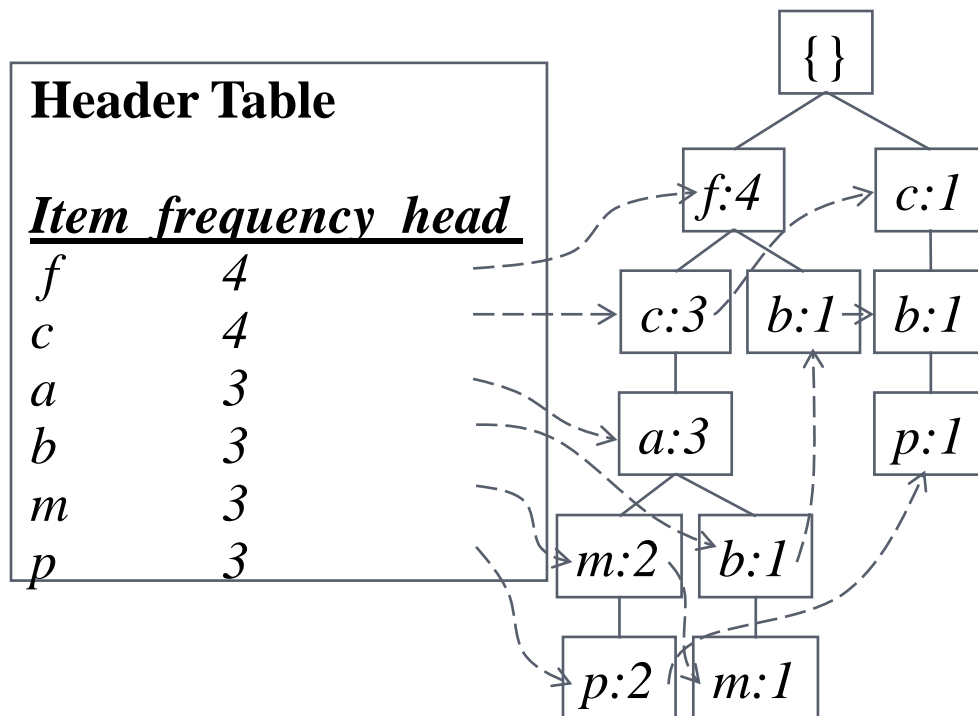
Patterns having c but no a nor b, m, p

Pattern f

- Completeness and non-redundance

4.2. FP-GROWTH EXAMPLE: FIND PATTERNS HAVING p FROM P-CONDITIONAL DATABASE

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



Conditional pattern bases

item *cond. pattern base*

c $f:3$

a $fc:3$

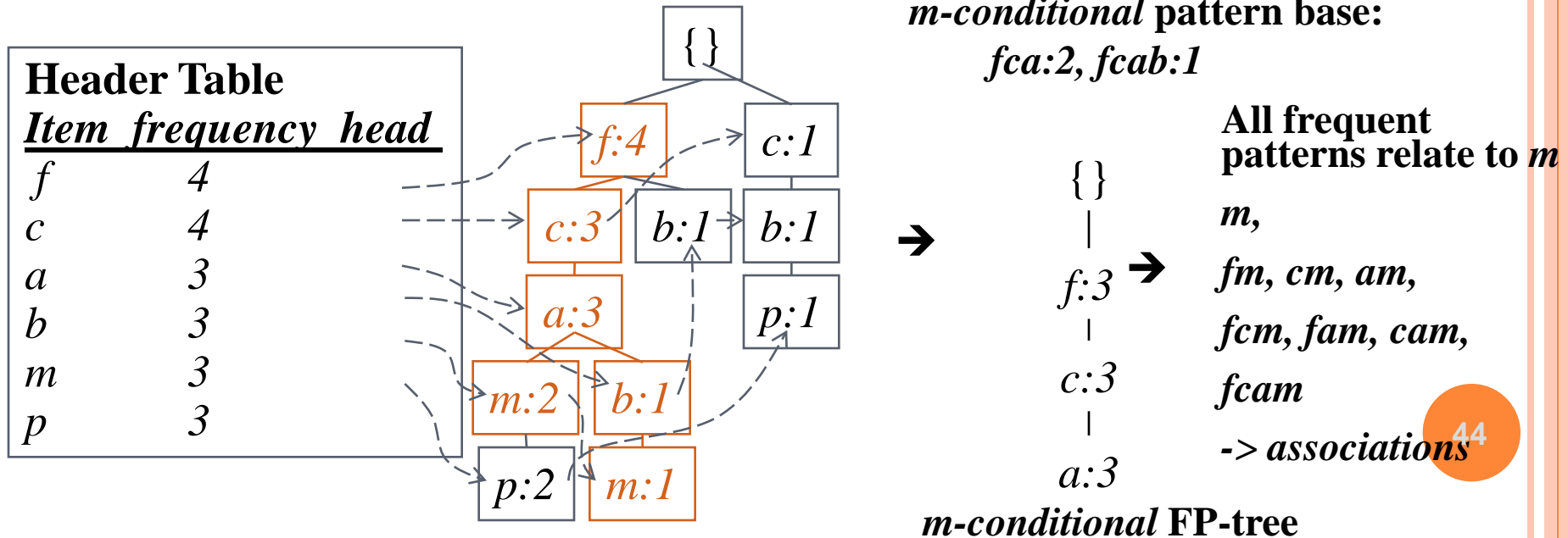
b $fca:1, f:1, c:1$

m $fca:2, fcab:1$

p $fcam:2, cb:1$

4.2. FP-GROWTH EXAMPLE: FROM CONDITIONAL PATTERN-BASES TO CONDITIONAL FP-TREES

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



4.2. FP-GROWTH EXAMPLE: RECURSION: MINING EACH CONDITIONAL FP-TREE

{ }
|
f:3
|
c:3
|
a:3

m-conditional FP-tree

Cond. pattern base of "am": (*fc:3*)

{ }
|
f:3
|
c:3

am-conditional FP-tree

Cond. pattern base of "cm": (*f:3*)

{ }
|
f:3

cm-conditional FP-tree

Cond. pattern base of "cam": (*f:3*)

{ }
|
f:3

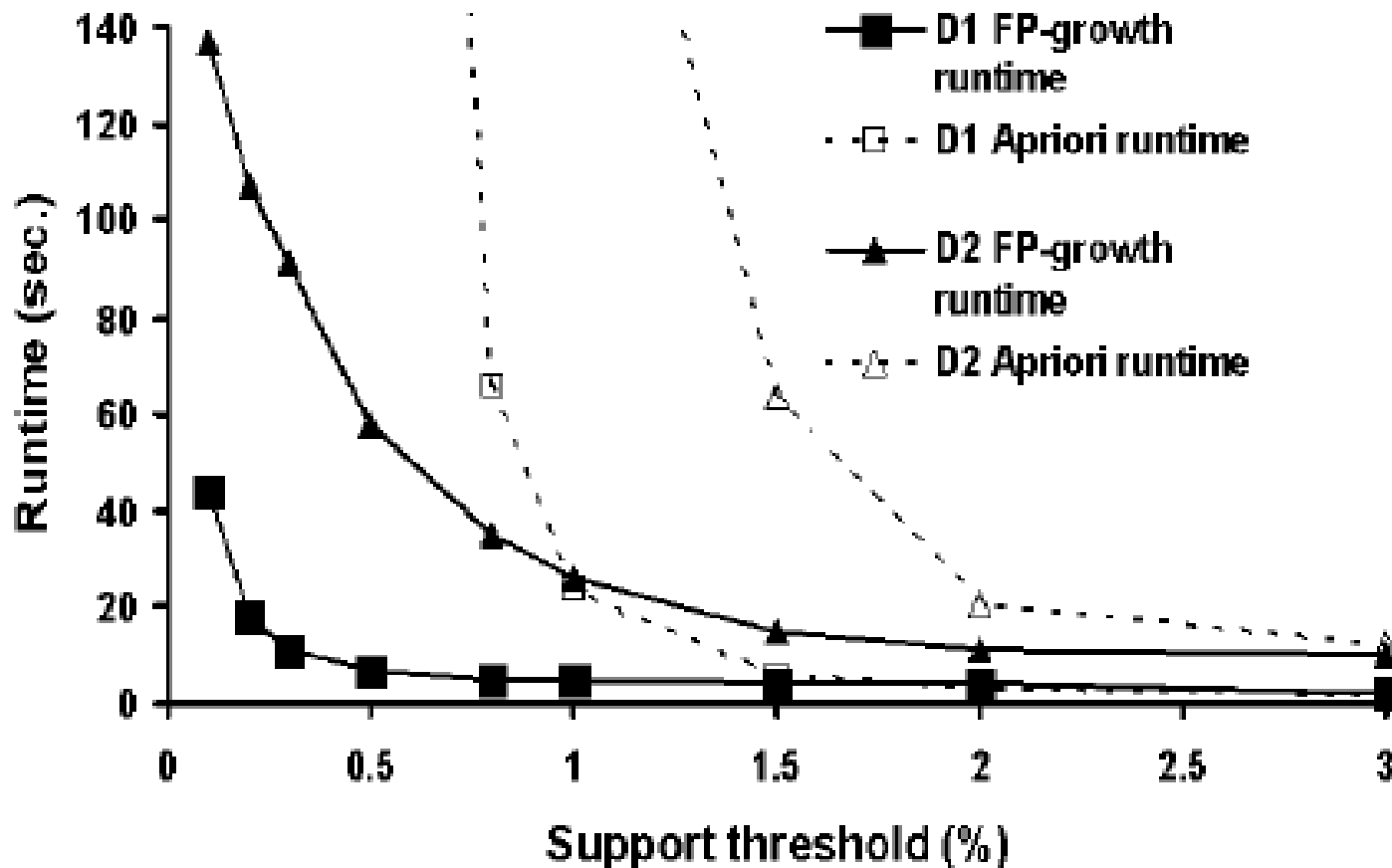
cam-conditional FP-tree

4.2. MINING FREQUENT ITEMSETS - FPGROWTH

- FP-Growth's main characteristics
 - Does not create candidate itemsets
 - FP-tree is a compact data structure
 - Reduce the cost for scanning the original dataset
 - Main costs come from building FP-tree (at the initial) and mining it
- Effective and scalable for mining long and sort frequent itemsets

4.2. MINING FREQUENT ITEMSETS - FPGROWTH

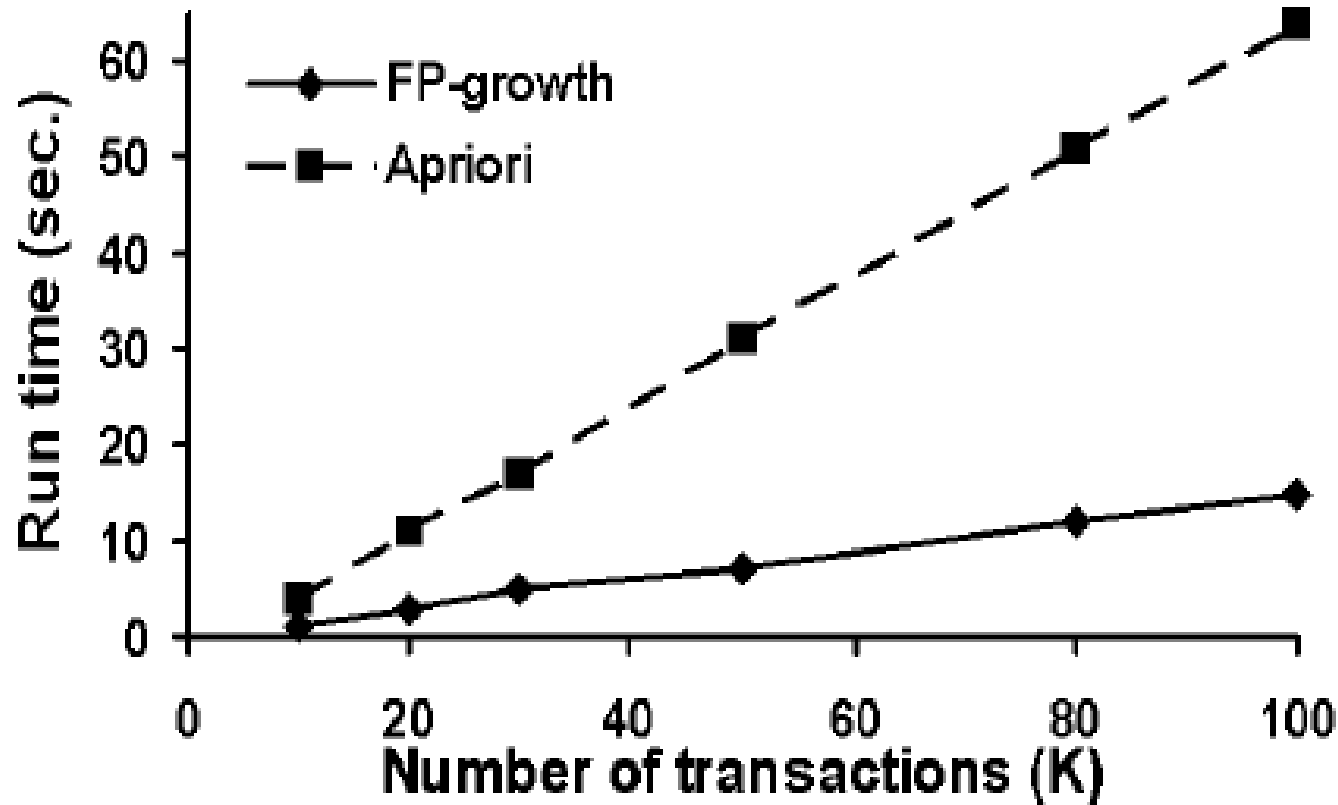
- Comparison between Apriori and FP-Growth



FP-growth scales well with support threshold

4.2. MINING FREQUENT ITEMSETS - FPGROWTH

- Comparison between Apriori and FP-Growth



Linearly scales with number of transactions

4.3 VARIANCE OF FREQUENT ITEMSET MINING

- Some variances of frequent itemsets
 - Frequent itemsets/subsequences/substructures
 - Closed frequent itemsets
 - Maximal frequent itemsets
 - Constrained frequent itemsets
 - Approximate frequent itemsets
 - Top-k frequent itemsets

5. MINING ASSOCIATION RULES FROM FREQUENT ITEMSETS

- Strong association rules $A \rightarrow B$
 - $\text{Support}(A \rightarrow B) = \text{Support}(A \cup B) \geq \text{min_sup}$
 - $\text{Confidence}(A \rightarrow B) = \text{Support}(A \cup B) / \text{Support}(A) = P(B | A) \geq \text{min_conf}$
- $\text{Support}(A \rightarrow B) = \text{Support_count}(A \cup B) \geq \text{min_sup}$
- $\text{Confidence}(A \rightarrow B) = P(B | A) = \text{Support_count}(A \cup B) / \text{Support_count}(A) \geq \text{min_conf}$

5. MINING ASSOCIATION RULES FROM FREQUENT ITEMSETS

- Process to generate strong association rules from frequent itemsets
 - For each frequent itemset l , create non-empty subsets of l .

$$\text{Support_count}(l) \geq \text{min_sup}$$

- For each non-empty subset s of l , create rule “ $s \rightarrow (l - s)$ ” if:

$$\text{Support_count}(l) / \text{Support_count}(s) \geq \text{min_conf}$$

5. MINING ASSOCIATION RULES FROM FREQUENT ITEMSETS

$$l = \{I1, I2, I5\}$$

nonempty subsets of l are $\{I1, I2\}$, $\{I1, I5\}$, $\{I2, I5\}$, $\{I1\}$, $\{I2\}$, and $\{I5\}$

$I1 \wedge I2 \Rightarrow I5$,	$confidence = 2/4 = 50\%$
$I1 \wedge I5 \Rightarrow I2$,	$confidence = 2/2 = 100\%$
$I2 \wedge I5 \Rightarrow I1$,	$confidence = 2/2 = 100\%$
$I1 \Rightarrow I2 \wedge I5$,	$confidence = 2/6 = 33\%$
$I2 \Rightarrow I1 \wedge I5$,	$confidence = 2/7 = 29\%$
$I5 \Rightarrow I1 \wedge I2$,	$confidence = 2/2 = 100\%$

Min_conf = 50%



{	$I1 \wedge I2 \Rightarrow I5$
	$I1 \wedge I5 \Rightarrow I2$
	$I2 \wedge I5 \Rightarrow I1$
	$I5 \Rightarrow I1 \wedge I2$

Frequent Patterns Generated

$\{I2, I5: 2\}$, $\{I1, I5: 2\}$, $\{I2, I1, I5: 2\}$
 $\{I2, I4: 2\}$
 $\{I2, I3: 4\}$, $\{I1, I3: 4\}$, $\{I2, I1, I3: 2\}$
 $\{I2, I1: 4\}$

6. MINING ASSOCIATION RULES BASED ON CONSTRAINTS

- Constraints
 - Instruct the processes of mining patterns and rules
 - Limit the search space of mining process
 - Some common constraints
 - Knowledge type constraints
 - Data constraints
 - Level/dimension constraints
 - Interestingness constraints
 - Rule constraints

6. MINING ASSOCIATION RULES BASED ON CONSTRAINTS

- Knowledge type constraints
 - What kind of knowledge we want to consider when mining rules

Ex. Mining association rules or correlation rules
- Data constraints
 - The characteristics of the data to be mined
- Level/dimension constraints
 - What kind of dimensions/features or level of abstractions that we want to know when mining association rules
- Interestingness constraints: Definition of the measure, what is the threshold,...
- Rule constraints: The types of rules to be mined

6. MINING ASSOCIATION RULES BASED ON CONSTRAINTS

- The mining process becomes more effective and efficient
 - Rules are mined based on business needs/user requirements -> More effective
 - Optimizers can be utilized to improve the efficiency

7. CORRELATION ANALYSIS

- Strong association rules $A \Rightarrow B$ are mined based on:
 - Occurrence frequency of A and B (*min_sup*)
 - The conditional probability of B based on A (*min_conf*)
 - *minsupport* and *minconfidence* are set in accordance with the subjective views of users
 - A large amount of rules will be returned.
 - Within 10,000 transactions, 6,000 ones for *computer games*, 7,500 for *videos*, and 4,000 for both *computer games* and *videos*
 - Buys(X, “*computer games*”) \Rightarrow Buys (X, “*videos*”)
[support = 40%, confidence = 66%]

7. CORRELATION ANALYSIS

- Correlation analysis for association rule $A \Rightarrow B$
 - Examine the correlation and dependency between A and B
 - Based on data statistics
 - Measures are objective, do not subjectively depend on users' views
- Withing 10,000 transactions, 6,000 for *computer games*, 7,500 for *videos*, and 4,000 for both *computer games* and *videos*
 - Buys(X, “*computer games*”) \Rightarrow Buys (X, “*videos*”) [support = 40%, confidence = 66%]
 - $P(\text{“}i\text{videos”}) = 75\% > 66\%$: “*computer games*” and “*videos*” are negatively correlated with each other.

7. CORRELATION ANALYSIS

- Correlation rules: $A \Rightarrow B$ [support, confidence, correlation]
 - correlation: measuring the correlation between A and B.
 - Common correlation measures: *lift*, χ^2 (Chi-square), *all_confidence*, *cosine*
 - *lift*: Validate the independent occurrence between A and B based on probability
 - χ^2 (Chi-square): Validate the independence between A and B based on expected value and the observed value.
 - *all_confidence*: Validate rules based on maximum support
 - *cosine*: similar to *lift* but it help to mitigate the dependency to the total number of transactions (i.e., the size of the dataset)
- *all_confidence* and *cosine* are good for large datasets, do not dependent on transactions that do not contain any itemsets being validated (null-transactions)

7. CORRELATION ANALYSIS

○ *lift*

- $lift(A, B) < 1$: A and B are negatively correlated
- $lift(A, B) > 1$: A and B are positively correlated
- $lift(A, B) = 1$: A and B are independent

$$lift(A, B) = \frac{P(A \cup B)}{P(A)P(B)} = P(B | A) / P(B) = confidence(A \Rightarrow B) / support(B)$$

	<i>game</i>	\overline{game}	Σ_{row}
<i>video</i>	4,000	3,500	7,500
\overline{video}	2,000	500	2,500
Σ_{col}	6,000	4,000	10,000

$$P(\{game\}) = 0.60,$$

$$P(\{video\}) = 0.75$$

$$P(\{game, video\}) = 0.40.$$

$$P(\{game, video\}) / (P(\{game\}) \times P(\{video\})) = 0.40 / (0.60 \times 0.75) = 0.89.$$

$lift(\{game\} \Rightarrow \{video\}) = 0.89 < 1 \rightarrow \{game\}$ and $\{video\}$ are negatively correlated.

8. SUMMARY

- Association rule mining
 - Considered as one of the most important contributions from database communities in KDD
- Rule types: logical/quantitative rules, single dimension/multiple dimensions, single level/multiple levels of abstraction, association/correlation rules
- Forms of item/pattern: Frequent itemsets/subsequences/substructures, Closed frequent itemsets, Maximal frequent itemsets, Constrained frequent itemsets, Approximate frequent itemsets, Top-k frequent itemsets
- Mine frequent itemsets: Apriori and FP-Growth algorithms

REFERENCES

- [1] Jiawei Han, Micheline Kamber, “Data Mining: Concepts and Techniques”, Second Edition, Morgan Kaufmann Publishers, 2006.
- [2] David Hand, Heikki Mannila, Padhraic Smyth, “Principles of Data Mining”, MIT Press, 2001.
- [3] David L. Olson, Dursun Delen, “Advanced Data Mining Techniques”, Springer-Verlag, 2008.
- [4] Graham J. Williams, Simeon J. Simoff, “Data Mining: Theory, Methodology, Techniques, and Applications”, Springer-Verlag, 2006.
- [5] ZhaoHui Tang, Jamie MacLennan, “Data Mining with SQL Server 2005”, Wiley Publishing, 2005.
- [6] Oracle, “Data Mining Concepts”, B28129-01, 2008.
- [7] Oracle, “Data Mining Application Developer’s Guide”, B28131-01, 2008.
- [8] Ian H. Witten, Eibe Frank, “Data mining : practical machine learning tools and techniques”, Second Edition, Elsevier Inc, 2005.
- [9] Florent Messegia, Pascal Poncelet & Maguelonne Teisseire, “Successes and new directions in data mining”, IGI Global, 2008.
- [10] Oded Maimon, Lior Rokach, “Data Mining and Knowledge Discovery Handbook”, Second Edition, Springer Science + Business Media, LLC 2005, 2010.

Q&A

quangtran@hcmut.edu.vn

2020/5/22

62