Soft Computing Association Analysis I

Dr. Chun-Hao Chen

Outline

1. Motivation & Problem Definition

2. Association Mining & Properties

3. Advanced Rule Mining Algorithm

4. Related Research Directions

Why Association Mining



✓ Extract useful information for managers making a better sales plan



Transactions from Supermarket

- ✓ What associations you can find from the following transaction
- ✓ Try to find them and figure out your reasons

TID	Items	
	(Dairyland chocolate milk), (Foremost chocolate milk),	
D_I	(Old Mills white bread),(Wonder white bread),	
	(Present chocolate cookies),(Linton green tea beverage).	
(Dairyland chocolate milk), (Foremost chocolate milk),		
D_2	(Dairyland plain milk), (Old Mills white bread),	
	(Wonder wheat bread), (Present lemon cookies), (77 lemon cookies).	
D	(Old Mills white bread), (Old Mills wheat bread),	
D_3	(77 chocolate cookies), (77 lemon cookies).	
D	(Dairyland chocolate milk), (Old Mills white bread),	
D_4	(77 chocolate cookies).	
D_5	(Old Mills white bread), (Wonder wheat bread).	
D_6	(Dairyland chocolate milk), (Foremost chocolate milk),	
	(Linton black tea beverage), (Nestle black tea beverage).	

Specific Associations



✓ e.g., {Dairyland chocolate milk, Old Mills white bread} is a specific association

TID	Items	
D_{I}	(Dairyland chocolate milk), (Foremost chocolate milk),	
	(Old Mills white bread), (Wonder white bread),	
	(Present chocolate cookies),(Linton green tea beverage).	
(Dairyland chocolate milk), (Foremost chocolate milk),		
D_2	(Dairyland plain milk), (Old Mills white bread),	
	(Wonder wheat bread), (Present lemon cookies), (77 lemon cookies).	
D	(Old Mills white bread), (Old Mills wheat bread),	
D_3	(77 chocolate cookies), (77 lemon cookies).	
	(Dairyland chocolate milk), (Old Mills white bread),	
D_4	(77 chocolate cookies).	
D_5	(Old Mills white bread), (Wonder wheat bread).	
D_6	(Dairyland chocolate milk), (Foremost chocolate milk),	
	(Linton black tea beverage), (Nestle black tea beverage).	

3 times out of 6

Generalized Associations

- ✓ Focus on the four items: milk, bread, cookies, beverage
- √ The six transactions can be transformed as follows

TID	Items
	(* * milk), (* * milk),
D_{I}	(* * <mark>bread</mark>),(* * bread),
	(* * cookies),(* * beverage).
	(* * milk), (* * milk),
D_2	(* * milk), (* * <mark>bread</mark>),
	(* * bread), (* * cookies), (* * cookies).
	(* * <mark>bread</mark>), (* * bread),
D_3	(* * cookies), (* * cookies).
D	(* * milk), (* * <mark>bread</mark>),
D_4	(* * cookies).
D_5	(* * bread), (* * bread).
D_6	(* * milk), (* * milk),
	(* * beverage), (* * beverage).

{bread, cookies} → 4 times out of 6

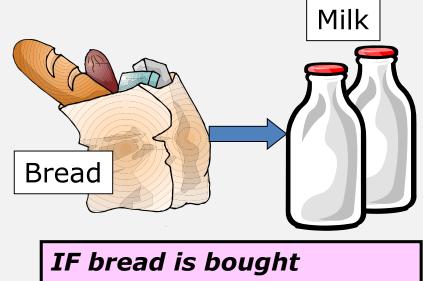
A Generalized Association

Association Rule Mining

✓ Proposed by Agrawal et al. (SIGMOD'93)

TID	Items
<i>T1</i>	milk, bread, cookies, beverage
<i>T</i> 2	milk, bread, cookies
<i>T3</i>	bread, cookies, beverage
<i>T4</i>	milk, bread
<i>T5</i>	milk, bread

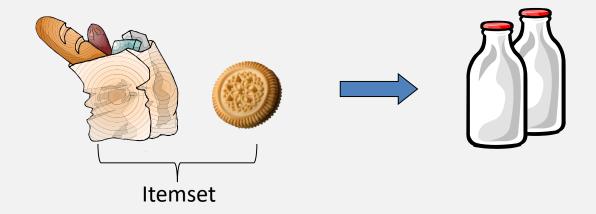




IF bread is bought then milk is bought

Problem Definition

- ✓ Given a set of transactions D
- ✓ Goal:
 - ullet To find association rules A ullet B from the dataset D that reach the minimum support lpha and minimum confidence λ constraints
 - Each of A and B is an itemsets (a set of items)



Apriori Algorithm

- ✓ The well-known algorithm & Proposed by Agrawal et al.
- √Three main phases
 - Phase 1: Define the two threshold minsup and minconf

- Phase 2: Find large itemsets
- Phase 3: Generate association rules

Example for Phase 2

Database		
TID	Items	
100	ACD	
200	ВСЕ	
300	ABCE	
400	ВЕ	

Scan Da<u>ta</u>base

C 1		
Itemset	Sup.	
{A}	2	
{B}	3	
{C}	3	
{D}	1	
{E}	3	

C 2	
Itemset	
{A B}	
{A C}	
{A E}	
{B C}	
{B E}	
{C E}	

Scan Da<u>ta</u>base

C 2		
Itemset	Sup.	
{A B}	1	
{A C}	2	
{A E}	1	
{B C}	2	
{B E}	3	
{C E}	2	



Scan Da<u>ta</u>base

C 3		
Itemset	Sup.	
{B C E}	2	

Large itemsets

L ₁		
Itemset	Sup	
{A}	2	
₹B}	3	
{C}	3	
{ E }	3	

L	2
Itemset	Sup.
{A C}	2
{B C}	2
{B E}	3
{C E}	2

L 3		
Itemset	Sup.	
{B C E}	2	

Example for Phase 3

Association rules	Confidence
IF BC THEN E	S(BCE)/S(BC)=2/2
IF BE THEN C	S(BCE)/S(BE)=2/3
IF CE THEN B	S(BCE)/S(CE)=2/2
IF B THEN CE	S(BCE)/S(B)=2/3
IF C THEN BE	S(BCE)/S(C)=2/3
IF E THEN BC	S(BCE)/S(E)=2/3
IF A THEN C	S(AC)/S(A)=2/2
IF C THEN A	S(AC)/S(C)=2/3
IF B THEN C	S(BC)/S(B)=2/3
IF C THEN B	S(BC)/S(C)=2/3
IF B THEN E	S(BE)/S(B)=3/3
IF E THEN B	S(BE)/S(E)=3/3
IF C THEN E	S(CE)/S(C)=2/3
IF E THEN C	S(CE)/S(E)=2/3

Step-by-Step

- ✓ Input Database (4 Transactions)
- ✓ Parameter setting
 - *■minsup* = 50%
 - *■ minconf* = 50%

Trans. D	atabase
TID	Items
100	ACD
200	ВСЕ
300	ABCE
400	ВЕ

Generate Large Itemsets

√ Generate Candidate 1-itemsets (C1)

Data	base
TID	Items
100	ACD
200	ВСЕ
300	ABCE
400	ВЕ

Scan Database

e.g. Item A appears in TID 100, 300

→	count =	2
----------	---------	---

C_{1}	
Itemset	Count
{A}	2
{B}	3
{C}	3
{D}	1
{E}	3

- ✓ Generate Large 1-itemsets (L1)
 - Check item's count larger than minimum support or not
 - •Total transactions = $4 \rightarrow \text{minimum count} = 2 (= 0.5*4)$

Database	
TID	Items
100	ACD
200	ВСЕ
300	ABCE
400	ВЕ

Scan Database

e.g. Item A appears in TID 100, 300

$$\rightarrow$$
 count = 2

C 1	
Itemset	Count
{A}	2
{B}	3
{C}	3
{D}	1
{E}	3

- ✓ Generate Large 1-itemsets (L₁)
 - Check item's count larger than minimum support or not
 - •Total transactions = $4 \rightarrow \text{minimum count} = 2 (= 0.5*4)$

C_1	
Itemset	Count
{A}	2
{B}	3
{C}	3
{D}	1
{E}	3

e.g.Count(A) = 2 ≥ 2Itemset A is L1

L	1
Itemset	Count
{A}	2
{B}	3
{C}	3
{E}	3

✓ Generate Candidate 2-itemsets (C₂) from large 1-itemsets (L₁)

L	1
Itemset	Count
{A}	2
{B}	3
{C}	3
{E}	3

e.g.

Large 1-itemset: A

 \rightarrow Generate three C_2

→ {AB}, {AC}, {AE}

C_2
Itemset
{A B}
{A C}
{A E}
{B C}
{B E}
{C E}

- ✓ Generate Large 2-itemsets (L₂)
 - Check itemset's count larger than minimum support or not
 - •minimum count = 2 = 0.5*4

C_2
Itemset
{A B}
{A C}
{A E}
{B C}
{B E}
{C E}

Scan Database

C_2	
Itemset	Count
{A B}	1
{A C}	2
{A E}	1
{B C}	2
{B E}	3
{C E}	2

e.g. Count(AB) = $1 \le 2$ → Delete {A B}

L ₂	
Itemset	Count
{A C}	2
{B C}	2
{B E}	3
{C E}	2

✓ Generate Candidate 3-itemsets (C₃) from L₂

L ₂	
Itemset	Count
{A C}	2
{B C}	2
{B E}	3
{C E}	2

e.g.
{BC} and {BE}

→generate C₃

→{BCE}

<i>C</i> ₃	
Itemset	
{B C E}	

- ✓ Generate Large 3-itemsets (L₃)
 - Check itemset's count larger than minimum support or not
 - •minimum count = 2 = 0.5*4

<i>C</i> ₃	
Itemset	
{B C E}	

Scan Database

Database	
TID	Items
100	ACD
200	ВСЕ
300	ABCE
400	ВЕ

C 3	
Itemset	Count
{B C E}	2

e.g.
Count(BCE) = $2 \le 2$
→ {BCE} is L3

L 3	
Itemset	Count
{B C E}	2

✓ All large itemsets

L ₁	
Itemset	Count
{A}	2
{B}	3
{C}	3
{E}	3

L ₂	
Itemset	Count
{A C}	2
{B C}	2
{B E}	3
{C E}	2

L 3	
Itemset Count	
{B C E}	2

Generate Association Rules

- ✓ Generate Candidate Association Rule
 - ✓ E.g. Large 2-itemset BC
 - √ Two candidate rules

1. B
$$\rightarrow$$
 C

2.
$$C \rightarrow B$$

L ₁		
Itemset	Count	
{A}	2	
{B}	3	
{C}	3	
{E}	3	

L ₂	
Itemset	Count
{A C}	2
{B C}	2
{B E}	3
{C E}	2

L 3	
Itemset	Count
{B C E}	2

Generate Association Rules (Cont.)

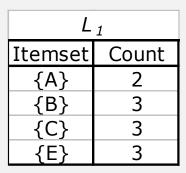
✓ Calculate confidence value of each rule

1. B
$$\rightarrow$$
 C

Conf(B
$$\rightarrow$$
 C) = P(BC)/P(B) = 2/3 = 0.66

2.
$$C \rightarrow B$$

Conf(C
$$\rightarrow$$
 B) = P(BC)/P(C) = 2/3 = 0.66



L_2		
Itemset	Count	
{A C}	2	
{B C}	2	
{B E}	3	
{C E}	2	

L 3		
Itemset	Count	
{B C E}	2	

Generate Association Rules (Cont.)

- ✓ Generate Association Rules
 - ●Minimum Confidence = 50%

1. B
$$\rightarrow$$
 C

Conf(B
$$\rightarrow$$
 C) = P(BC)/P(B) = 2/3 = 0.66



Rule₁: B \rightarrow C, sup=0.5, conf=0.66

2. C
$$\rightarrow$$
 B

Conf(C
$$\rightarrow$$
 B) = P(BC)/P(C) = 2/3 = 0.66



Rule₂: C \rightarrow B, sup=0.5, conf=0.66

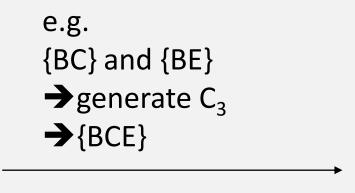
Generate Association Rules (Cont.)

✓ After generating all association rules

Association rules	Confidence
IF BC THEN E	S(BCE)/S(BC)=2/2
IF BE THEN C	S(BCE)/S(BE)=2/3
IF CE THEN B	<i>S(BCE)/S(CE)=2/2</i>
IF B THEN CE	S(BCE)/S(B)=2/3
IF C THEN BE	S(BCE)/S(C)=2/3
IF E THEN BC	S(BCE)/S(E)=2/3
IF A THEN C	S(AC)/S(A)=2/2
IF C THEN A	S(AC)/S(C)=2/3
IF B THEN C	S(BC)/S(B)=2/3
IF C THEN B	S(BC)/S(C)=2/3
IF B THEN E	S(BE)/S(B)=3/3
IF E THEN B	S(BE)/S(E)=3/3
IF C THEN E	S(CE)/S(C)=2/3
IF E THEN C	S(CE)/S(E)=2/3

Discussion

L 2	
Itemset	Sup.
{A C}	2
{B C}	2
{B E}	3
{C E}	2

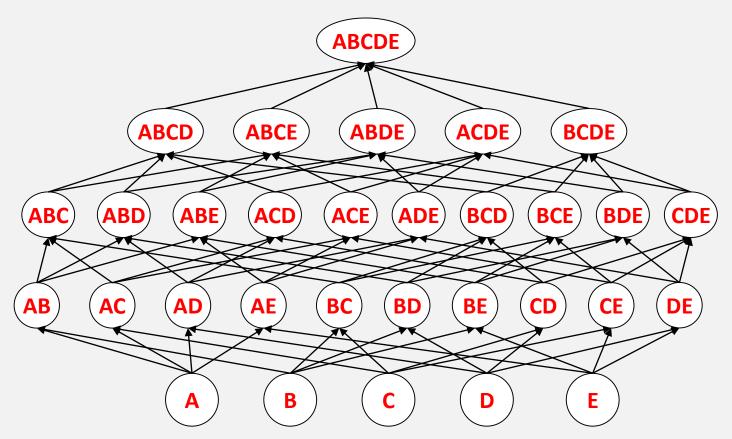


Itemset

Why generate only a candidate 3-itemset?

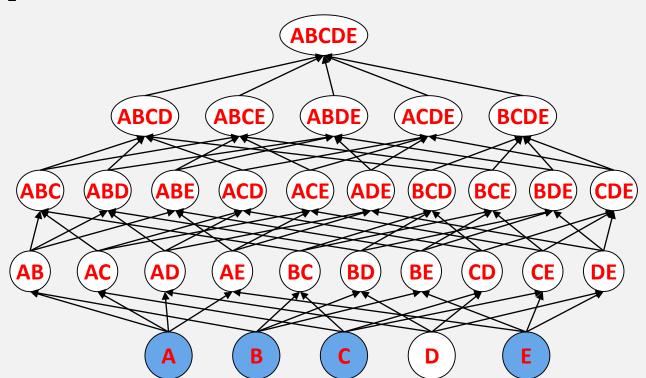
Downward Closure Property

✓ Subsets of a large k-itemset are also frequent.





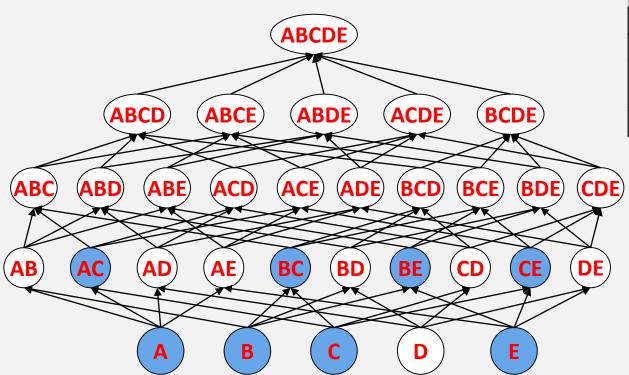
- ✓ Continue Previous Example
 - \bullet L₁: A, B, C and E



L ₁	
Itemset	Count
{A}	2
{B}	3
{C}	3
{E}	3

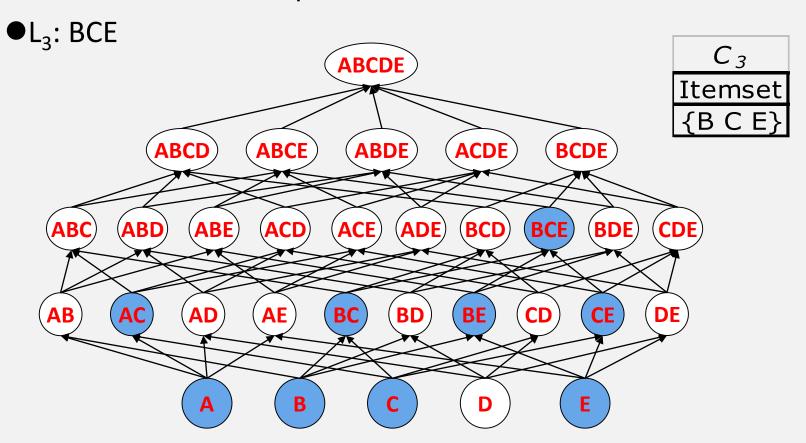
✓ Continue Previous Example



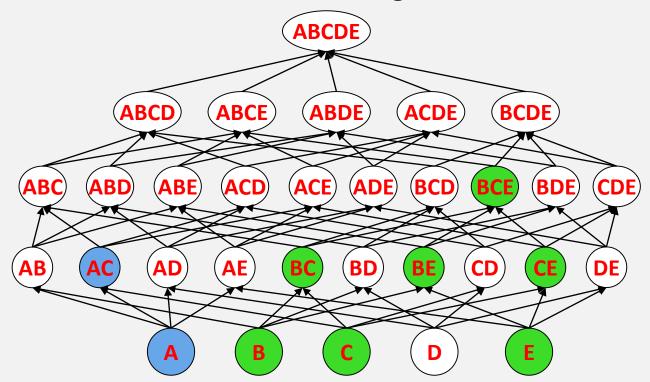


L ₂	
Itemset	Count
{A C}	2
{B C}	2
{B E}	3
{C E}	2

✓ Continue Previous Example

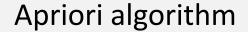


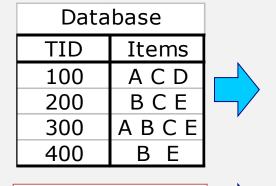
- ✓ BCE is a large 3-itemset
 - ●Then, BC, BE, CE, B, C, E are all large itemsets



Flowchart of Association Rule Mining







Min. Sup.

Min. Conf.

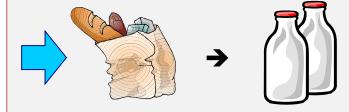


Candidate Itemsets



Large Itemsets

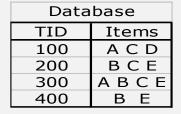
Association Rule



Discussion

Do you find any problem of Apriori Algorithm?

Problem 1: Scan Database



Scan Database

C_1	
Itemset	Sup.
{A}	2
{B}	3
{C}	3
{D}	1
{E}	3

C_2
Itemset
{A B}
{A C}
{A E}
{B C}
{B E}
{C E}



C 2	
Itemset	Sup.
{A B}	1
{A C}	2
{A E}	1
{B C}	2
{B E}	3
{C E}	2



Scan Database

C 3		
Itemset	Sup.	
{B C E}	2	

Large itemsets

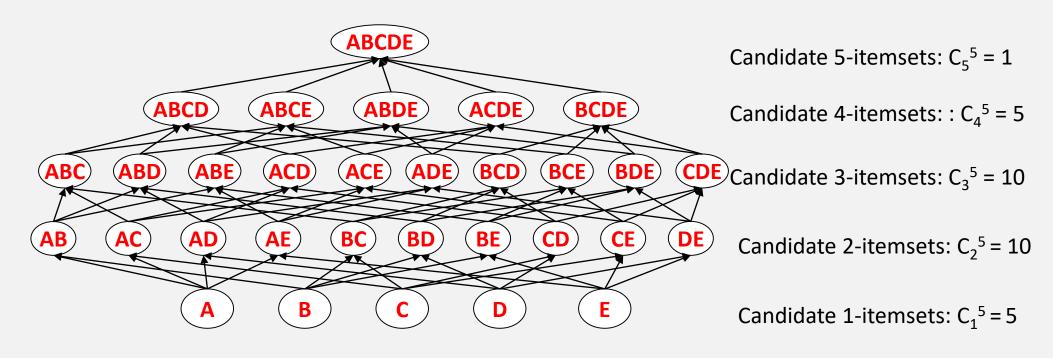
L 1		
Itemset	Sup.	
{ A }	2	
₹ B }	3	
{C}	3	
{F}	3	

L 2		
Itemset	Sup.	
{A C}	2	
{B C}	2	
{B E}	3	
{C E}	2	

L	L 3		
Itemset	Sup.		
{B C E}	2		

Problem 2: Candidate Itemsets





$$C_n^k = \binom{n}{k} = \frac{P_n^k}{k!} = \frac{n!}{k!(n-k)!}$$



$$C_2^{100} = 100*99/1*2 = 50*99$$
 $C_2^{1000} = 1000*999/1*2 = 500*999$
 $C_2^{10000} = 10000*9999/1*2 = 5000*9999$

Main Problems of Apriori Algorithm

- ✓ Scan the database several times
- ✓ Time-consuming for generating candidate itemsets

Could We Avoid Those Problems?

- √ Solution
 - Proposed by Han et al.
 - Jiawei Han, Jian Pei and Yiwen Yin, "Mining Frequent Patterns without Candidate Generation," SIGMOD Conference, pp. 1-12, 2000



FP-Growth Algorithm

Flowchart of the FP-growth Algorithm



Transactional Database		
Transaction ID	Items Purchased	
10	A, B, C	
20	A, C	
30	A, D	
40	B. E. F	



Candidate 1-itemsets



min_sup

Frequent 1-itemsets



Sort L1 in descending frequency



Construct FP-Tree

All Frequent Itemsets



FP-Growth Algorithm - Step-by-Step

- ✓ Two Phases
 - Phase 1: Constructing FP-tree
 - Phase 2: Executing FP-Growth mining method
 - ✓ Mining from FP-tree

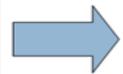
Phase 1: Construction of FP-Tree

- ✓ Three steps
 - Scan DB to find L1
 - Sort L1 in descending frequency
 - Scan DB again
 - ✓ Construct FP-tree

Step 1

- Scanning DB to find L1
- Example: min_sup=60%*5=3 trans

TID	Item bought
100	a, c, d, f, g, i, m, p
200	a, b, c, f, l, m, o
300	b, f, h, j, o, w
400	b, c, k, p, s
500	a, c, e, f, l, m, n, p



 $L1 = \{a, b, c, f, m, p\}$

item	frequency	item	frequency
а	3	j	1
b	3	k	1
С	4	1	2
d	1	m	3
е	1	0	2
f	4	р	3
g	1	s	1
h	1	w	1
i	1		

Step 2



Sort L1 in descending frequency

$$min_sup = 3$$

- L1 = {a:3, b:3, c:4, f:4, m:3, p:3}L1' = {f:4, c:4, a:3, b:3, m:3, p:3}
- Sort DB

TID	Item bought
100	a, c, d, f, g, i, m, p
200	a, b, c, f, l, m, o
300	b, f, h, j, o, w
400	b, c, k, p, s
500	a, c, e, f, l, m, n, p



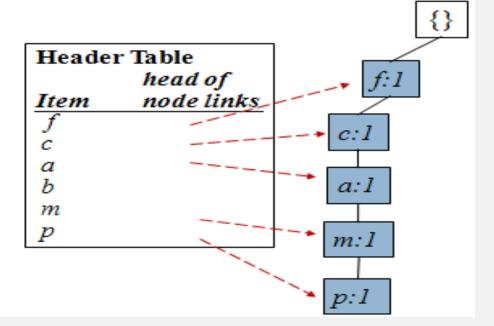
TID	Frequent items
100	f, c, a, m, p
200	f, c, a, b, m
300	f, b
400	c, b, p
500	f, c, a, m, p

Step 3

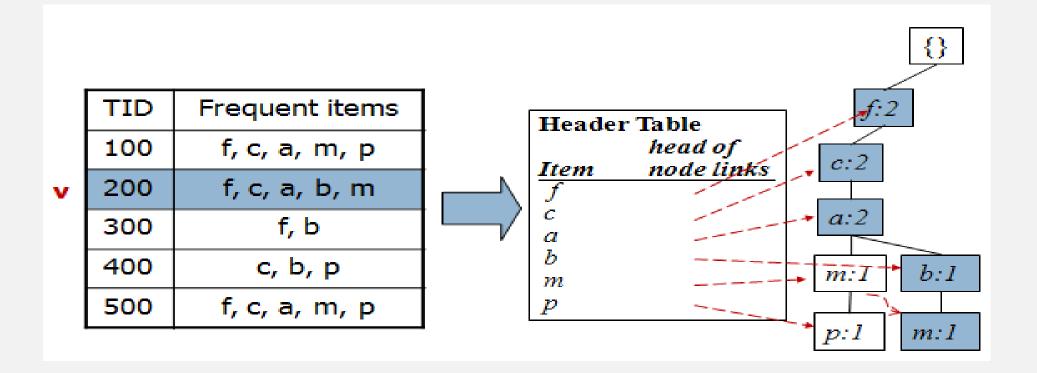
- Scanning DB again
 - Construct FP-tree

TID	Frequent items
100	f, c, a, m, p
200	f, c, a, b, m
300	f, b
400	c, b, p
500	f, c, a, m, p





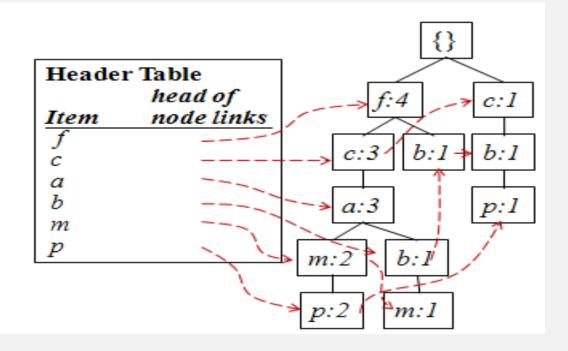
Step 3 (Cont.)



Step 3 (Cont.)

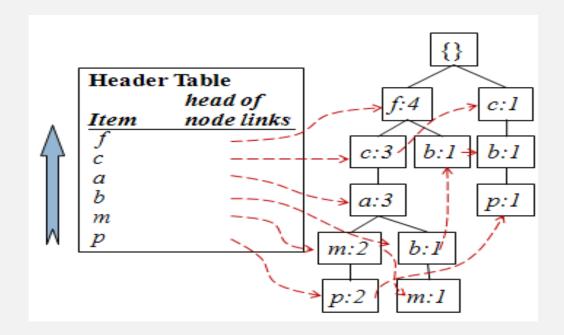


TID	Frequent items
100	f, c, a, m, p
200	f, c, a, b, m
300	f, b
400	c, b, p
500	f, c, a, m, p



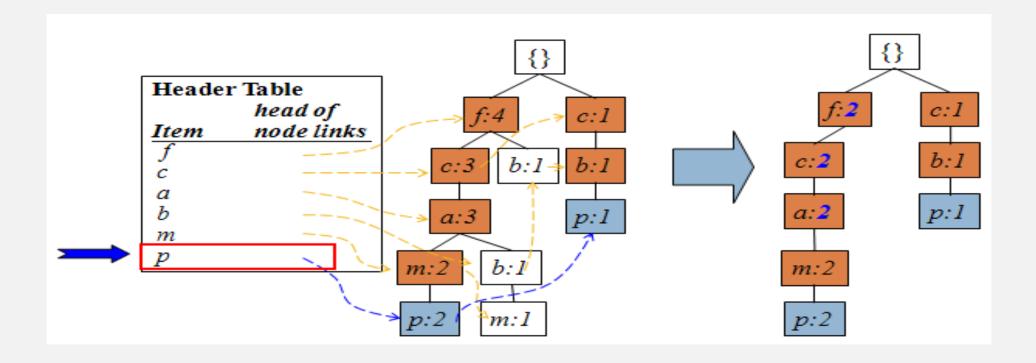
Phase 2: FP-Growth

- ✓ Mining frequent patterns from FP-tree
- ✓ Processing frequent items
 - One by one
 - Bottom up



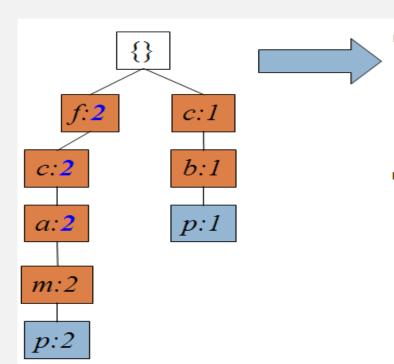
- ✓ Each item
 - Generating its own conditional FP-tree

Example for *p*



Example for p (Cont.)



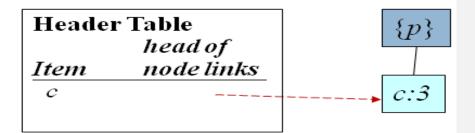


P's conditional pattern base

$$\rightarrow$$
 {(c:3, f:2, a:2, m:2, b:1)}

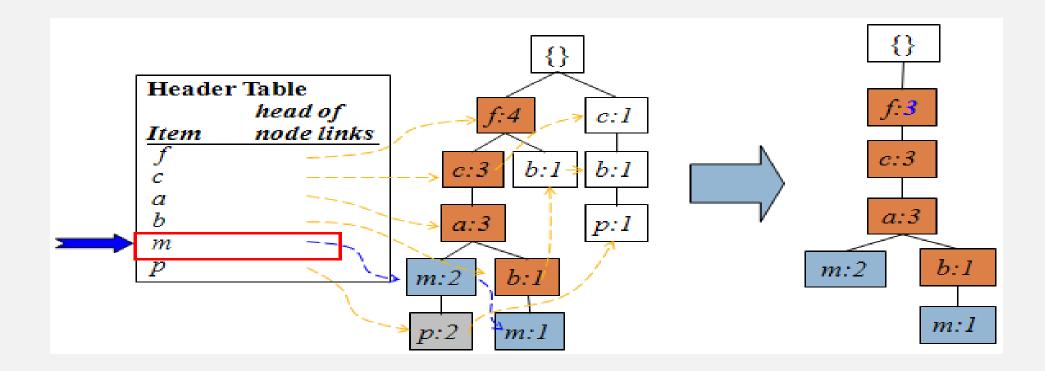
min sup = 3

Constructing p's conditional FP-tree:



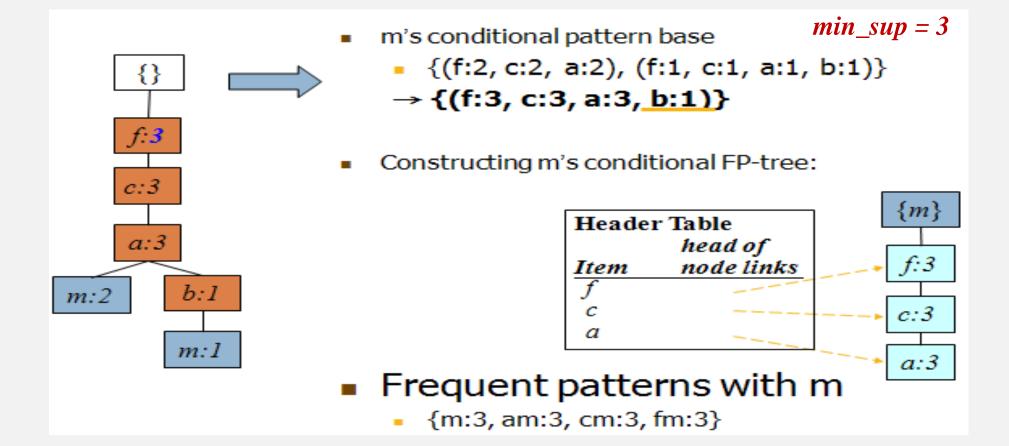
- Frequent patterns with p
 - {p:3, cp:3}

Example for *m*



Example for m (Cont.)



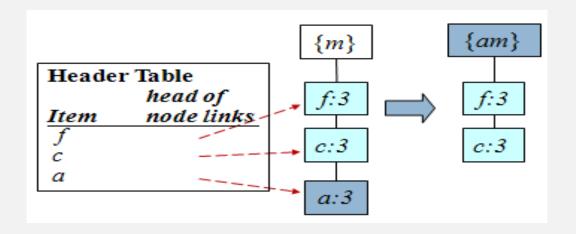


Example for am

- ✓ Frequent patterns with *m*
 - •{m:3, am:3, cm:3, fm:3}
- ✓ Recursively constructing conditional FP-tree
 - •am, cm, fm
 - Bottom up
- √am first
 - Prefix

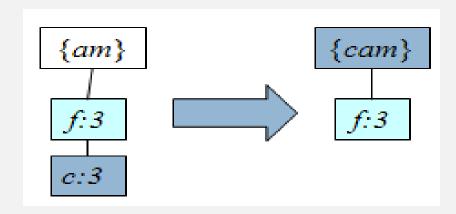
- Large items
 - √f:3, c:3
- Conditional FP tree for am
- •Large itemsets with am

```
✓ cam:3, fam:3
```



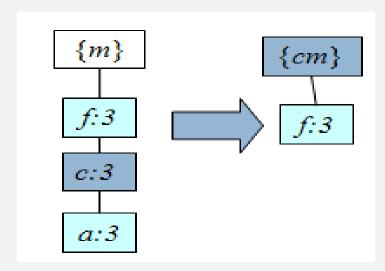
Example for *cam* and *fam*

- ✓ Frequent patterns with *am*
 - •{am:3, cam:3, fam:3}
- ✓ Recursively constructing conditional FP-tree
 - •cam, fam
 - Bottom up
- ✓ cam first
 - •fcam:3
- √fam next
 - ●null



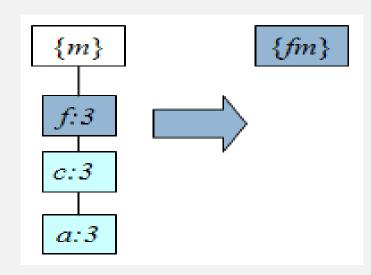
Example for *cm*

- ✓ Frequent patterns with *m*
 - •{m:3, am:3, cm:3, fm:3}
- √cm
 - Prefix
 - **√**(f:3)
 - Large items
 - √f:3
 - Conditional FP tree for cm
 - •Large itemsets with *cm*
 - √fcm:3



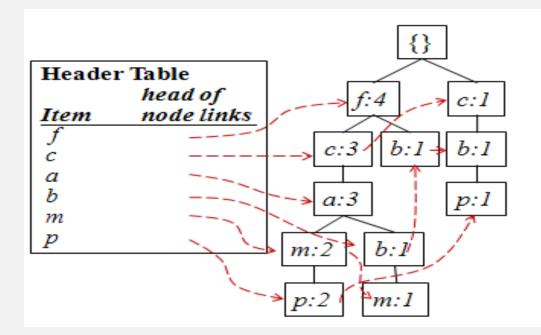
Example for fm

- ✓ Frequent patterns with *m*
 - •{m:3, am:3, cm:3, fm:3}
- √fm
 - Null



- ✓ All large itemsets with *m*
 - •m:3, fm:3, cm:3, am:3, fcm:3, fam:3, cam:3, fcam:3

Result



Item	frequent itemsets
f	f:4
C	c:4, fc:3
а	a:3, ca:3, fa:3, fca:3
b	b:3
m	m:3, fm:3, cm:3, am:3, fcm:3, fam:3, cam:3, fcam:3
р	p:3, cp:3

Summary

- ✓ A Novel data structure → Frequent-Patten tree (FP-tree)
 - Compressing database to store frequent patterns
 - Saving costly scanning database
- √FP-growth mining method
 - Efficient mining frequent patterns from FP-tree
- ✓ Without candidate generation
 - Faster than the Apriori algorithm
- ✓ Scan database only twice

Related Research Direction



Could you extend the association rule to other types of association rule?

Problems of Association Rule

- ✓ In real world applications
 - Quantitative value may exist in transaction



TID	Items
T1	(milk, 6); (bread, 4); (cookies, 7), (beverage, 7).
<i>T</i> 2	(milk, 7); (bread, 7); (cookies, 12).
<i>T3</i>	(bread, 8); (cookies, 12); (beverage, 6).
<i>T4</i>	(milk, 2); (bread, 3).
<i>T5</i>	(milk, 3); (bread, 8).
—	

Quantitative value

First problem

→ How to handle Quantitative Transactions?

Quantitative Association Rules

- ✓ Proposed by Srikant *et al.* (SIGMOD'96)
 - Extend from Apriori Algorithm

TID	Items
<i>T1</i>	(Age: 50, Car, 5)
<i>T</i> 2	(Age: 24, Car, 1)
<i>T3</i>	(Age: 28, Car, 1)
T4	(Age: 22, Car, 0)
<i>T5</i>	(Age: 35, Car, 4)



TID	Items
<i>T1</i>	(Age: [50, 59], Car: [3, 5])
<i>T</i> 2	(Age: [20, 29], Car: [0, 1])
<i>T3</i>	(Age: [20, 29], Car: [0, 1])
<i>T4</i>	(Age: [20, 29], Car: [0, 1])
<i>T5</i>	(Age: [30, 39], Car: [3, 5])

IF Age is [20,... 29], then Number of Car is [0, 1]. Sup. = 60%, Conf. = 66.6%



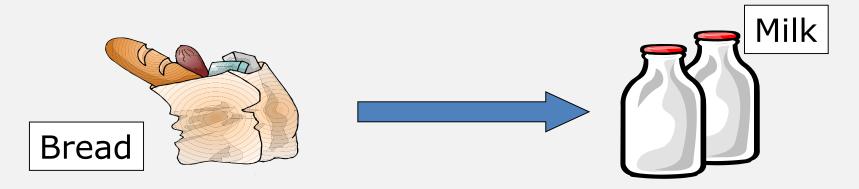
Age: 25



Car

Fuzzy Association Rules

- ✓ Proposed by Kuok et al. (SIGMOD Record, 1998)
 - Another approach to handle quantitative value



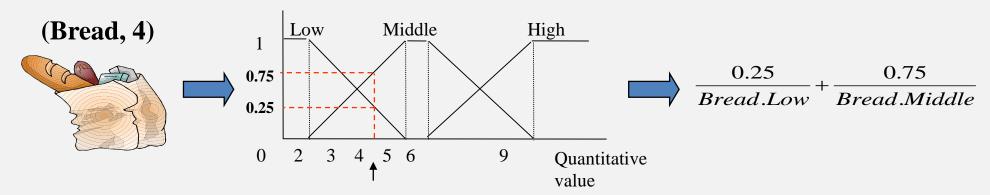
IF middle amount of bread is bought
Then high amount of milk is bought

Linguistic Term

Fuzzy Association Rules (Cont.)

- √ How to handle quantitative data?
- √e.g.

Membership functions of bread



Benefits of Fuzzy Association Rules

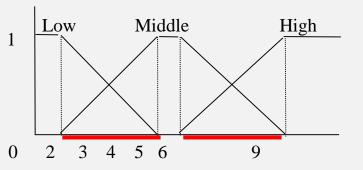
- √Three Advantages (Kuok et al.)
 - ✓ Understandable to Human
 - √ Handling Quantitative Value Well (Shape boundary problem)
 - ✓ Deriving Extra Information

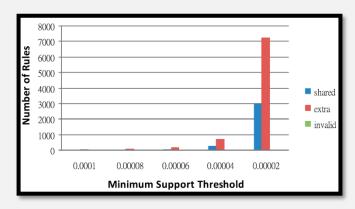
(Bread, 4)



Bread.Middle Bread.High

Membership functions





Problems of Association Rule (Cont.)

- ✓ Second problem
 - A minimum support for all item is not appropriate

√e.g.

Supermarket



Price: 1,000

Sales: 50

Gain: 50,000

Pan

Milk

Price: 20

Sales: 200

Gain: 4,000

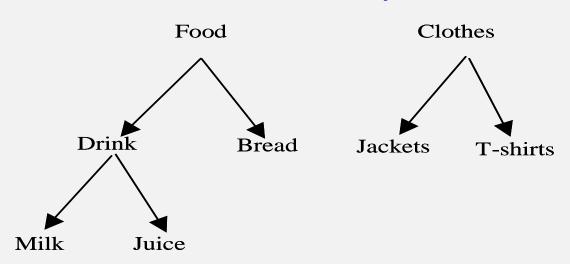
High Utility Item

Different Items Different Minimum Supports

Problems of Association Rule (Cont.)

- ✓ Only the terminal items can appear in transaction data
- √ However, items may have taxonomy (as follows)

Taxonomy



Conclusion

- ✓ Association Rule Mining algorithm
 - Binary association rule (Apriori)
 - An efficient rule mining algorithm (FP-growth)
 - Property of the Apriori algorithm
 - Related problems of Apriori algorithm and possible solutions