# DATA MINING 1 Pattern Mining & Association Rule Mining

Dino Pedreschi, Riccardo Guidotti

Revisited slides from Lecture Notes for Chapter 5 "Introduction to Data Mining", 2nd Edition by Tan, Steinbach, Karpatne, Kumar



# **Association Rules - Module Outline**

- What are association rules (AR) and what are they used for:
  - The paradigmatic application: Market Basket Analysis
  - The single dimensional AR (intra-attribute)

# **Market Basket Analysis: The Context**

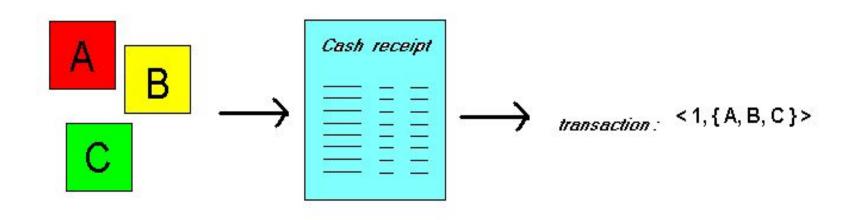
 Analyzing customer purchasing habits by finding associations and correlations between the different items that customers place in their "shopping basket"



# **Market Basket Analysis: The Context**

• Given: a database of customer transactions, where each transaction is a set of items.

Goal: Find groups of items which are frequently purchased together.



# **Goal of MBA**

- Extract information on purchasing behavior
- Actionable information: can suggest
  - new store layouts
  - new product assortments
  - which products to put on promotion
- MBA applicable whenever a customer purchases multiple things in proximity
  - credit cards
  - services of telecommunication companies
  - banking services
  - medical treatments

# MBA: applicable to many other contexts

### Telecommunication:

Each customer is a transaction containing the set of customer's phone calls

### Atmospheric phenomena:

Each time interval (e.g. a day) is a transaction containing the set of observed event (rains, wind, etc.)

Etc.

# **Association Rules**

- Express how product/services relate to each other, and tend to group together
- "if a customer purchases three-way calling, then will also purchase call-waiting"
- simple to understand
- actionable information: bundle three-way calling and call-waiting in a single package
- Examples.
  - Rule form: "Body → Head [support, confidence]".
  - buys(x, "diapers") → buys(x, "beers") [0.5%, 60%]
  - major(x, "CS") and takes(x, "DB") → grade(x, "A") [1%, 75%]

# Useful, trivial, unexplicable

- Useful: "On Thursdays, grocery store consumers often purchase diapers and beer together".
- Trivial: "Customers who purchase maintenance agreements are very likely to purchase large appliances".
- Unexplicable: "When a new hardaware store opens, one of the most sold items is toilet rings."

# **Apriori**

# **Association Rule Mining**

 Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction

### Market-Basket transactions

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

**Example of Association Rules** 

```
{Diaper} \rightarrow {Beer},

{Milk, Bread} \rightarrow {Eggs,Coke},

{Beer, Bread} \rightarrow {Milk},
```

Implication means co-occurrence, not causality!

# **Definition: Frequent Itemset**

### Itemset

- A collection of one or more items
  - Example: {Milk, Bread, Diaper}
- k-itemset
  - An itemset that contains k items

### Support count (σ)

- Frequency of occurrence of an itemset
- E.g.  $\sigma(\{Milk, Bread, Diaper\}) = 2$

### Support

- Fraction of transactions that contain an itemset
- E.g. s({Milk, Bread, Diaper}) = 2/5

### Frequent Itemset

• An itemset whose support is greater than or equal to a minsup threshold

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

# **Definition: Association Rule**

### Association Rule

- An implication expression of the form
   X → Y, where X and Y are itemsets
- Example:{Milk, Diaper} → {Beer}

Pula	Evalu	uation	$N/I_{c}$	strice
Ruie	-vall	Janon	IVIE	#1111C/S

- Support (s)
  - Fraction of transactions that contain both X and Y
- Confidence (c)
  - Measures how often items in Y appear in transactions that contain X

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

### Example:

$$\{Milk, Diaper\} \Rightarrow \{Beer\}$$

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

# **Association rules - module outline**

- How to compute AR
  - Basic Apriori Algorithm and its optimizations
  - Multi-Dimension AR (inter-attribute)
  - Quantitative AR

# **Association Rule Mining Task**

- Given a set of transactions T, the goal of association rule mining is to find all rules having
  - support ≥ *minsup* threshold
  - confidence ≥ minconf threshold
- Brute-force approach:
  - List all possible association rules
  - Compute the support and confidence for each rule
  - Prune rules that fail the minsup and minconf thresholds
  - ⇒ Computationally prohibitive!

# **Mining Association Rules**

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

### Example of Rules:

```
 \begin{aligned} &\{\text{Milk,Diaper}\} \rightarrow \{\text{Beer}\} \text{ } (\text{s=0.4, c=0.67}) \\ &\{\text{Milk,Beer}\} \rightarrow \{\text{Diaper}\} \text{ } (\text{s=0.4, c=1.0}) \\ &\{\text{Diaper,Beer}\} \rightarrow \{\text{Milk}\} \text{ } (\text{s=0.4, c=0.67}) \\ &\{\text{Beer}\} \rightarrow \{\text{Milk,Diaper}\} \text{ } (\text{s=0.4, c=0.67}) \\ &\{\text{Diaper}\} \rightarrow \{\text{Milk,Beer}\} \text{ } (\text{s=0.4, c=0.5}) \\ &\{\text{Milk}\} \rightarrow \{\text{Diaper,Beer}\} \text{ } (\text{s=0.4, c=0.5}) \end{aligned}
```

### **Observations:**

- All the above rules are binary partitions of the same itemset: {Milk, Diaper, Beer}
- Rules originating from the same itemset have identical support but can have different confidence
- Thus, we may decouple the support and confidence requirements

# **Mining Association Rules**

- Two-step approach:
  - 1. Frequent Itemset Generation
    - Generate all itemsets whose support ≥ minsup
  - 2. Rule Generation
    - Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

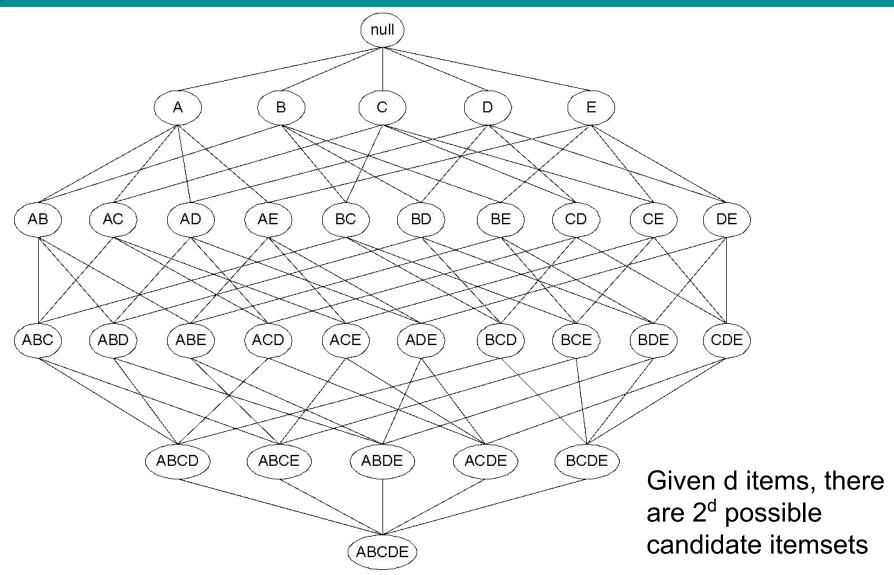
Frequent itemset generation is still computationally expensive

# **Basic Apriori Algorithm**

### **Problem Decomposition**

- 1. Find the *frequent itemsets*: the sets of items that satisfy the support constraint
  - A subset of a frequent itemset is also a frequent itemset, i.e., if {A,B} is a frequent itemset, both {A} and {B} should be a frequent itemset
  - Iteratively find frequent itemsets with cardinality from 1 to k (k-itemset)
- 2. Use the frequent itemsets to generate association rules.

# **Frequent Itemset Generation**



# **Frequent Itemset Generation**

- Brute-force approach:
  - Each itemset in the lattice is a candidate frequent itemset
  - Count the support of each candidate by scanning the database

    Transactions

    List of

			<b>Candidates</b>	
	TID	Items		
	1	Bread, Milk		·
T	2	Bread, Diaper, Beer, Eggs		ı
Ň	3	Milk, Diaper, Beer, Coke		M
1	4	Bread, Milk, Diaper, Beer		1
$\downarrow$	5	Bread, Milk, Diaper, Coke		V

- Match each transaction against every candidate
- Complexity ~ O(NMw) => Expensive since M = 2<sup>d</sup> !!!

# **Frequent Itemset Generation Strategies**

- Reduce the number of candidates (M)
  - Complete search: M=2<sup>d</sup>
  - Use pruning techniques to reduce M
- Reduce the number of transactions (N)
  - Reduce size of N as the size of itemset increases
- Reduce the number of comparisons (NM)
  - Use efficient data structures to store the candidates or transactions
  - No need to match every candidate against every transaction

# **Reducing Number of Candidates**

- Apriori principle:
  - If an itemset is frequent, then all of its subsets must also be frequent
- Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \ge s(Y)$$

- Support of an itemset never exceeds the support of its subsets
- This is known as the anti-monotone property of support

**Illustrating Apriori Principle** null С Ε D ВС BD BE CE DE Found to be Infrequent BDE ABC (ABD) ACD ACE ADE BCD BCE CDE ABCD ABCE ABDE ACDE BCDE **Pruned** ABCDE supersets

TID	Items
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Coke, Diaper, Milk



Items (1-itemsets)

Item	Count	
Bread	4	
Coke	2	
Milk	4	
Beer	3	
Diaper	4	
Eggs	1	

Minimum Support = 3

TID	Items
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Coke, Diaper, Milk



Items (1-itemsets)

Item	Count	
Bread	4	
Coke	2	
Milk	4	
Beer	3	
Diaper	4	
Eggs	1	

Minimum Support = 3

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset
{Bread,Milk}
{Bread, Beer }
{Bread,Diaper}
{Beer, Milk}
{Diaper, Milk}
{Beer,Diaper}

Minimum Support = 3

TID	Items
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Coke, Diaper, Milk

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Beer, Bread}	2
{Bread,Diaper}	3
{Beer,Milk}	2
{Diaper,Milk}	3
{Beer,Diaper}	3

1 Bread, Milk
2 Beer, Bread, Diaper, Eggs
3 Beer, Coke, Diaper, Milk
4 Beer, Bread, Diaper, Milk
5 Bread, Coke, Diaper, Milk

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Beer, Bread}	2
{Bread,Diaper}	3
{Beer,Milk}	2
{Diaper,Milk}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

TID

Items

Bread, Milk

Beer, Bread, Diaper, Eggs

Beer, Coke, Diaper, Milk

Beer, Bread, Diaper, Milk

Bread, Coke, Diaper, Milk

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3



Triplets (3-itemsets)

Itemset
{ Beer, Diaper, Milk}
{ Beer, Bread, Diaper}
{Bread, Diaper, Milk}
{ Beer, Bread, Milk}

(No need to generate candidates involving Bread, Beer or Milk, Beer)

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Beer, Bread}	2
{Bread,Diaper}	3
{Beer,Milk}	2
{Diaper,Milk}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

TID

Items

Bread, Milk

Beer, Bread, Diaper, Eggs

Beer, Coke, Diaper, Milk

Beer, Bread, Diaper, Milk

Bread, Coke, Diaper, Milk

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3



Triplets (3-itemsets)

Itemset	Count
{ Beer, Diaper, Milk}	2
{ Beer,Bread, Diaper}	2
{Bread, Diaper, Milk}	2
{Beer, Bread, Milk}	1

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Beer, Bread}	2
{Bread,Diaper}	3
{Beer,Milk}	2
{Diaper,Milk}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

TID

Items

Bread, Milk

Beer, Bread, Diaper, Eggs

Beer, Coke, Diaper, Milk

Beer, Bread, Diaper, Milk

Bread, Coke, Diaper, Milk

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3



Triplets (3-itemsets)

Itemset	Count
{ Beer, Diaper, Milk}	2
{ Beer,Bread, Diaper}	2
{Bread, Diaper, Milk}	2
{Beer, Bread, Milk}	1

# **Apriori Algorithm**

- F<sub>k</sub>: frequent k-itemsets
- L<sub>k</sub>: candidate k-itemsets
- Algorithm
  - Let k=1
  - Generate F<sub>1</sub> = {frequent 1-itemsets}
  - Repeat until F<sub>k</sub> is empty
    - Candidate Generation: Generate L<sub>k+1</sub> from F<sub>k</sub>
    - Candidate Pruning: Prune candidate itemsets in  $L_{k+1}$  containing subsets of length k that are infrequent
    - Support Counting: Count the support of each candidate in  $L_{k+1}$  by scanning the DB
    - Candidate Elimination: Eliminate candidates in  $L_{k+1}$  that are infrequent, leaving only those that are frequent =>  $F_{k+1}$

# Candidate Generation: $F_{k-1} \times F_{k-1}$ Method

- Merge two frequent (k-1)-itemsets if they have a common element (if k = 2)
- Merge two frequent (k-1)-itemsets if their first (k-2) items are identical (if k > 2)
- F<sub>3</sub> = {ABC,ABD,ABE,ACD,BCD,BDE,CDE}
  - Merge( $\underline{AB}C$ ,  $\underline{AB}D$ ) =  $\underline{AB}CD$
  - Merge( $\underline{AB}C$ ,  $\underline{AB}E$ ) =  $\underline{AB}CE$
  - Merge( $\underline{AB}D$ ,  $\underline{AB}E$ ) =  $\underline{AB}DE$
  - Do not merge(<u>ABD</u>,<u>ACD</u>) because they share only prefix of length 1 instead of length 2

# **Candidate Pruning**

- Let F<sub>3</sub> = {ABC,ABD,ABE,ACD,BCD,BDE,CDE} be the set of frequent 3-itemsets
- L<sub>4</sub> = {ABCD,ABCE,ABDE} is the set of candidate 4-itemsets generated (from previous slide)
- Candidate pruning
  - Prune ABCE because ACE and BCE are infrequent
  - Prune ABDE because ADE is infrequent
- After candidate pruning:  $L_{\Delta} = \{ABCD\}$

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3
,	11

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3



Triplets (3-itemsets)



Use of  $F_{k-1}xF_{k-1}$  method for candidate generation results in only one 3-itemset. This is eliminated after the support counting step.

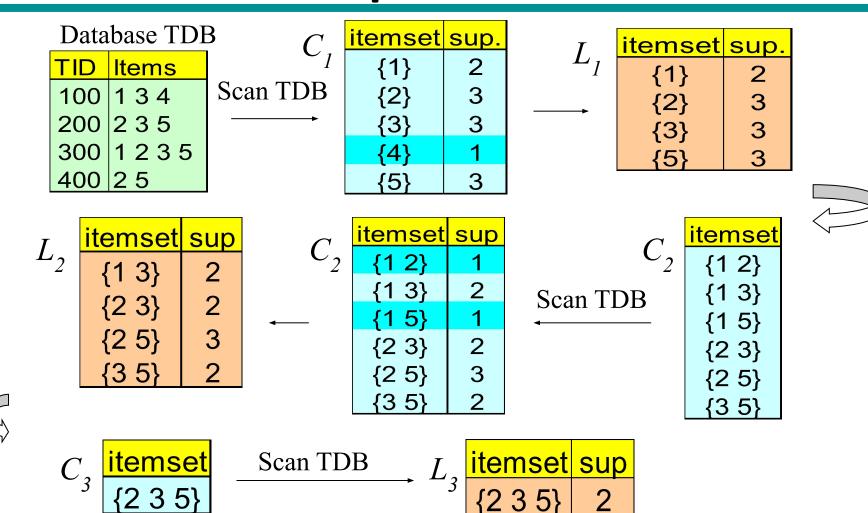
# **Support Counting of Candidate Itemsets**

- Scan the database of transactions to determine the support of each candidate itemset
- Must match every candidate itemset against every transaction, which is an expensive operation

TID	Items
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Coke, Diaper, Milk



# Apriori Execution Example (min\_sup = 2)



# **Rule Generation**

- Given a frequent itemset L, find all non-empty subsets  $f \subset L$  such that  $f \to L f$  satisfies the minimum confidence requirement
  - If {A,B,C,D} is a frequent itemset, candidate rules:

ABC 
$$\rightarrow$$
D, ABD  $\rightarrow$ C, ACD  $\rightarrow$ B, BCD  $\rightarrow$ A, A  $\rightarrow$ BCD, B  $\rightarrow$ ACD, C  $\rightarrow$ ABD, D  $\rightarrow$ ABC AB  $\rightarrow$ CD, AC  $\rightarrow$  BD, AD  $\rightarrow$  BC, BC  $\rightarrow$ AD, BD  $\rightarrow$ AC, CD  $\rightarrow$ AB,

• If |L| = k, then there are  $2^k - 2$  candidate association rules (ignoring  $L \to \emptyset$  and  $\emptyset \to L$ )

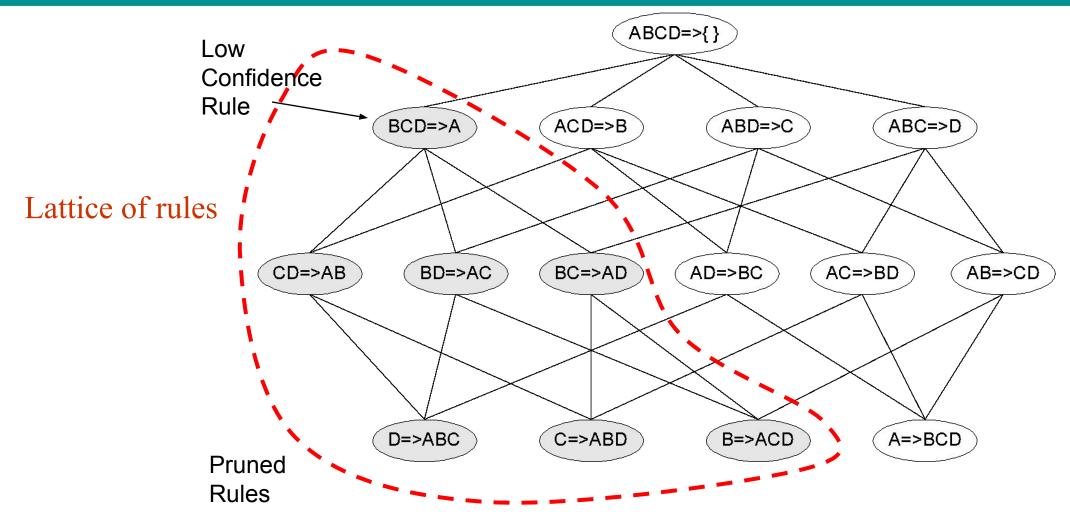
## **Rule Generation**

- In general, confidence does not have an anti-monotone property c(ABC →D) can be larger or smaller than c(AB →D)
- But confidence of rules generated from the same itemset has an anti-monotone property
  - E.g., Suppose {A,B,C,D} is a frequent 4-itemset:

$$c(ABC \rightarrow D) \ge c(AB \rightarrow CD) \ge c(A \rightarrow BCD)$$
  
 $s(A,B,C,D) / s(A,B,C) s(A,B,C,D) / s(A,B) s(A,B,C,D) / s(A)$   
 $s(A,B,C) \le s(A,B) \le s(A)$ 

- Confidence is anti-monotone w.r.t. number of items on the RHS of the rule
- If  $c(ABC \rightarrow D) < min\_conf$  it is useless to calculate  $c(AB \rightarrow CD)$

# Rule Generation for Apriori Algorithm



# **Maximal Frequent Itemset**

An itemset is maximal frequent if it is frequent and none of its immediate supersets is frequent null Maximal В C D **Itemsets** ΑE ВС BD BE CD CE DE AC ABD ABE ACD ACE ADE BCD BCE BDE CDE ABCD ABCE ABDE ACDE BCDE Infrequent **Itemsets** Border ABCD

### **Closed Itemset**

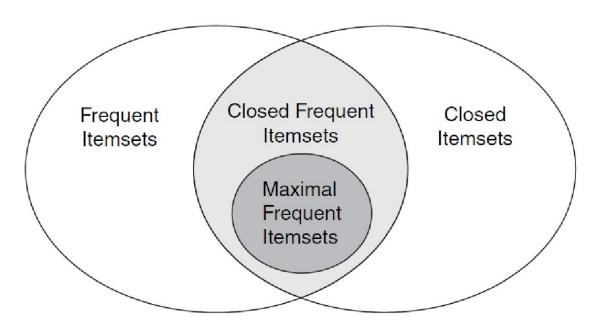
- An itemset X is closed if none of its immediate supersets has the same support as the itemset X.
- X is not closed if at least one of its immediate supersets has support count as X.

TID	Items
1	{A,B}
2	{B,C,D}
3	$\{A,B,C,D\}$
4	{A,B,D}
5	$\{A,B,C,D\}$

Itemset	Support
{A}	4
{B}	5
{C}	3
{D}	4
{A,B}	4
{A,C}	2
{A,D}	3
{B,C}	3
{B,D}	4
{C,D}	3

Itemset	Support
$\{A,B,C\}$	2
$\{A,B,D\}$	3
$\{A,C,D\}$	2
$\{B,C,D\}$	2
$\{A,B,C,D\}$	2

## **Maximal vs Closed Itemsets**



**Figure 5.18.** Relationships among frequent, closed, closed frequent, and maximal frequent itemsets.

## **Pattern Evaluation**

Association rule algorithms can produce large number of rules

- Interestingness measures can be used to prune/rank the patterns
- In the original formulation, support & confidence are the only measures used

## **Computing Interestingness Measure**

• Given  $X \rightarrow Y$  or  $\{X,Y\}$ , information needed to compute interestingness can be obtained from a contingency table

#### Contingency table

	Y	Y	
X	f <sub>11</sub>	f <sub>10</sub>	f <sub>1+</sub>
X	f <sub>01</sub>	f <sub>00</sub>	f <sub>o+</sub>
	f <sub>+1</sub>	f <sub>+0</sub>	N

f<sub>11</sub>: support of X and Y  $f_{10}$ : support of X and  $\overline{Y}$ 

 $f_{01}$ : support of  $\overline{X}$  and  $\overline{Y}$   $f_{00}$ : support of  $\overline{X}$  and  $\overline{Y}$ 

#### Used to define various measures

support, confidence, Gini, entropy, etc.

## **Drawback of Confidence**

Custo mers	Tea	Coffee	
C1	0	1	
C2	1	0	
C3	1	1	•••
C4	1	0	

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

Association Rule: Tea → Coffee

Confidence  $\cong$  P(Coffee|Tea) = 150/200 = 0.75

Confidence > 50%, meaning people who drink tea are more likely to drink coffee than not drink coffee

So, this rule seems reasonable

## **Drawback of Confidence**

	Coffee	Coffee	
Tea	150	50	200
Tea	650	150	800
	800	200	1000

Association Rule: Tea → Coffee

Confidence= P(Coffee|Tea) = 150/200 = 0.75

but P(Coffee) = 0.8, which means knowing that a person drinks tea reduces the probability that the person drinks coffee!

 $\rightarrow$  Note that P(Coffee | Tea) = 650/800 = 0.8125

## Measure for Association Rules

- So, what kind of rules do we really want?
  - Confidence( $X \rightarrow Y$ ) should be sufficiently high
    - To ensure that people who buy X will more likely buy Y than not buy Y
  - Confidence( $X \rightarrow Y$ ) > support(Y)
    - Otherwise, rule will be misleading because having item X actually reduces the chance of having item Y in the same transaction
    - Is there any measure that capture this constraint?
      - Answer: Yes. There are many of them.

# Statistical Relationship between X and Y

 The criterion confidence(X → Y) = support(Y)

#### is equivalent to:

- $\bullet P(Y \mid X) = P(Y)$
- $P(X,Y) = P(X) \times P(Y)$  (X and Y are independent)

If  $P(X,Y) > P(X) \times P(Y) : X \& Y$  are positively correlated

If  $P(X,Y) < P(X) \times P(Y) : X & Y are negatively correlated$ 

#### Measures that take into account statistical dependence

$$Lift = \frac{P(Y \mid X)}{P(Y)}$$

$$Interest = \frac{P(X,Y)}{P(X)P(Y)}$$

$$PS = P(X,Y) - P(X)P(Y)$$

$$\phi - coefficient = \frac{P(X,Y) - P(X)P(Y)}{\sqrt{P(X)[1 - P(X)]P(Y)[1 - P(Y)]}}$$

# **Example: Lift/Interest**

	Coffee	Coffee	
Tea	150	50	200
Tea	650	150	800
	800	200	1000

Association Rule: Tea → Coffee

Confidence= P(Coffee|Tea) = 0.75but P(Coffee) = 0.8

 $\Rightarrow$  Interest = 0.15 / (0.2×0.8) = 0.9375 (< 1, therefore is negatively associated)

## **Continuous and Categorical Attributes**

How to apply association analysis to non-asymmetric binary variables?

Gender	Gender $\cdots$ Age		Annual	No of hours spent	No of email	Privacy
			Income	online per week	accounts	Concern
Female	5.5.5	26	90K	20	4	Yes
Male	1.1.1	51	135K	10	2	No
Male		29	80K	10	3	Yes
Female		45	120K	15	3	Yes
Female		31	95K	20	5	Yes
Male	*:*:	25	55K	25	5	Yes
Male		37	100K	10	1	No
Male		41	65K	8	2	No
Female	77.	26	85K	12	1	No
	***	***				

Example of Association Rule:

 $\{Gender=Male, Age ∈ [21,30)\} \rightarrow \{No of hours online ≥ 10\}$ 

# **Handling Categorical Attributes**

Example: Internet Usage Data

Gender	Level of Education	State	Computer at Home	Online Auction	Chat Online	Online Banking	Privacy Concerns
Female	Graduate	Illinois	Yes	Yes	Daily	Yes	Yes
Male	College	California	No	No	Never	No	No
Male	Graduate	Michigan	Yes	Yes	Monthly	Yes	Yes
Female	College	Virginia	No	Yes	Never	Yes	Yes
Female	Graduate	California	Yes	No	Never	No	Yes
Male	College	Minnesota	Yes	Yes	Weekly	Yes	Yes
Male	College	Alaska	Yes	Yes	Daily	Yes	No
Male	High School	Oregon	Yes	No	Never	No	No
Female	Graduate	Texas	No	No	Monthly	No	No
	***		***		3.65	***	

```
{Level of Education=Graduate, Online Banking=Yes}

→ {Privacy Concerns = Yes}
```

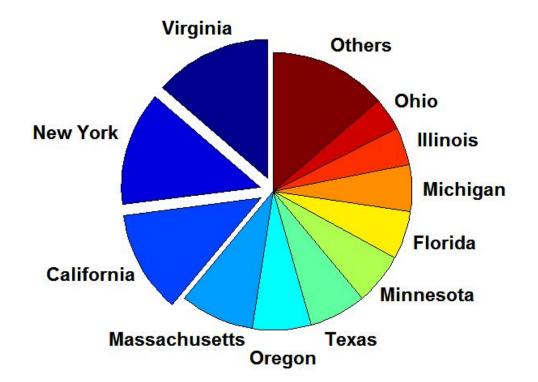
# **Handling Categorical Attributes**

• Introduce a new "item" for each distinct attribute-value pair

Male	Female	Education	Education	Education		Privacy	Privacy
		= Graduate	= College	= High School		= Yes	= No
0	1	1	0	0		1	0
1	0	0	1	0	5.55	0	1
1	0	1	0	0	***	1	0
0	1	0	1	0		1	0
0	1	1	0	0	* * *	1	0
1	0	0	1	0		1	0
1	0	0	0	0		0	1
1	0	0	0	1	****	0	1
0	1	1	0	0	***	0	1
		***	***	***	4.4.5	***	

# **Handling Categorical Attributes**

- Some attributes can have many possible values
  - Many of their attribute values have very low support
    - Potential solution: Aggregate the low-support attribute values

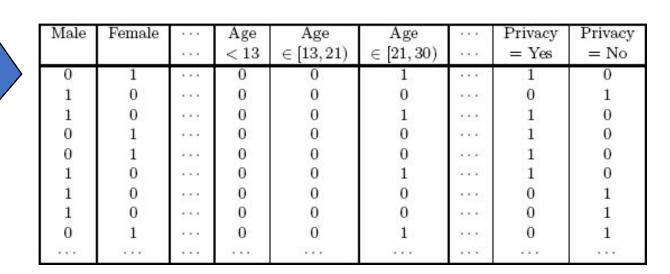


# **Handling Continuous Attributes**

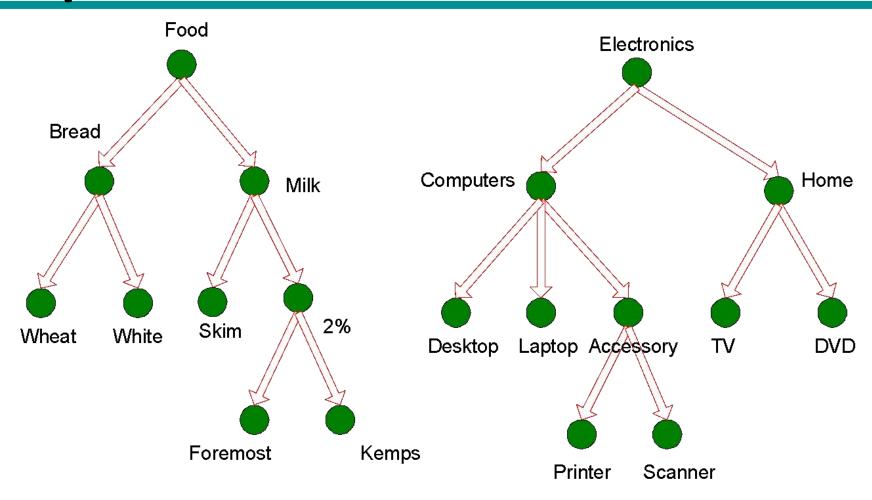
- Different methods:
  - Discretization-based
  - Statistics-based
  - Non-discretization based
    - minApriori
- Different kinds of rules can be produced:
  - {Age  $\in$  [21,30), No of hours online  $\in$  [10,20)}  $\rightarrow$  {Chat Online = Yes}
  - {Age  $\subseteq$  [21,30), Chat Online = Yes}  $\rightarrow$  No of hours online:  $\mu$ =14,  $\sigma$ =4

## **Discretization-based Methods**

Gender	***	Age	Annual Income	No of hours spent online per week	No of email accounts	Privacy Concern
Female	100	26	90K	20	4	Yes
Male	2.23	51	135K	10	2	No
Male		29	80K	10	3	Yes
Female		45	120K	15	3	Yes
Female	***	31	95K	20	5	Yes
Male	XXX	25	55K	25	5	Yes
Male	100	37	100K	10	1	No
Male		41	65K	8	2	No
Female		26	85K	12	1	No
		2.55			***	



# **Concept Hierarchies**



## **Multi-level Association Rules**

- Why should we incorporate concept hierarchy?
  - Rules at lower levels may not have enough support to appear in any frequent itemsets
  - Rules at lower levels of the hierarchy are overly specific
    - e.g., skim milk → white bread, 2% milk → wheat bread, skim milk → wheat bread, etc.
       are indicative of association between milk and bread
  - Rules at higher level of hierarchy may be too generic

## **Multi-level Association Rules**

 Approach 1: Extend current association rule formulation by augmenting each transaction with higher level items

Original Transaction: {skim milk, wheat bread}

Augmented Transaction: {skim milk, wheat bread, milk, bread, food}

#### Issues:

- Items that reside at higher levels have much higher support counts
- if support threshold is low, too many frequent patterns involving items from the higher levels
- Increased dimensionality of the data

## **Multi-level Association Rules**

- Approach 2:
  - Generate frequent patterns at highest level first
  - Then, generate frequent patterns at the next highest level, and so on
- Issues:
  - I/O requirements will increase dramatically because we need to perform more passes over the data
  - May miss some potentially interesting cross-level association patterns

# Is Apriori Fast Enough?

- The core of the Apriori algorithm:
  - Use frequent (k-1)-itemsets to generate candidate frequent k-itemsets
  - Use database scan and pattern matching to collect counts for the candidate itemsets
- The bottleneck of Apriori: candidate generation
  - Huge candidate sets:
    - 10<sup>4</sup> frequent 1-itemset will generate 10<sup>7</sup> candidate 2-itemsets
    - To discover a frequent pattern of size 100, e.g.,  $\{a_1, a_2, ..., a_{100}\}$ , one needs to generate  $2^{100} \approx 10^{30}$  candidates.
  - Multiple scans of database:
    - Needs (n +1) scans, n is the length of the longest pattern

# **FP-Growth**

## Mining Frequent Patterns Without Candidate Generation

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

TID	Items bought
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	$\{a, f, c, e, l, p, m, n\}$

 $min\_support = 3$ 

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

Hea	der Table
Iten	frequency head
f	4
c	4
a	3
b	3
m	3
p	3

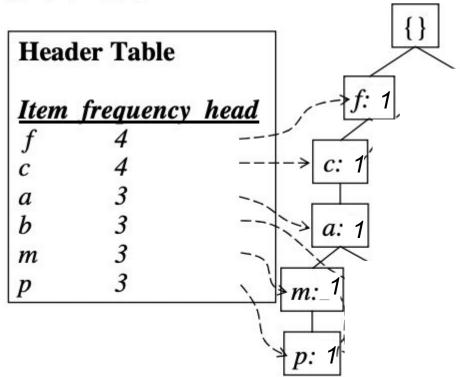
TID	Items bought	(ordered) frequent items	
100	$\{f, a, c, d, g, i, m, p\}$	$\{f, c, a, m, p\}$	$min_support = 3$
200	$\{a, b, c, f, l, m, o\}$	$\{f, c, a, b, m\}$	- 11
300	$\{b, f, h, j, o\}$	$\{f, b\}$	
400	$\{b, c, k, s, p\}$	$\{c, b, p\}$	
500	$\{a, f, c, e, l, p, m, n\}$	$\{f, c, a, m, p\}$	

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

Head	ler Table	
<u>Item</u>	frequency	head
f	4	_
c	4	
a	3	
b	3	
m	3	
p	3	

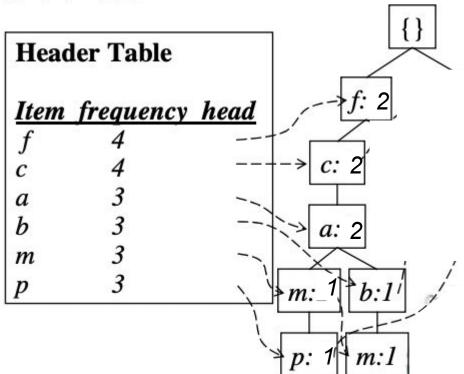
<b>TID</b>	Items bought	(ordered) frequent items	
100	$\{f, a, c, d, g, i, m, p\}$	$\{f, c, a, m, p\}$	$min_support = 3$
200	$\{a, b, c, f, l, m, o\}$	$\{f, c, a, b, m\}$	
300	$\{b, f, h, j, o\}$	$\{f, b\}$	
400	$\{b, c, k, s, p\}$	$\{c, b, p\}$	
500	$\{a, f, c, e, l, p, m, n\}$	$\{f, c, a, m, p\}$	

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



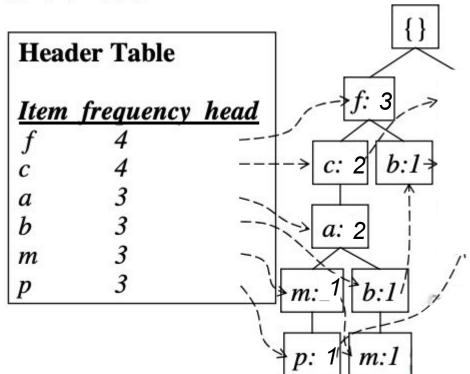
<b>TID</b>	Items bought	(ordered) frequent items	
100	$\{f, a, c, d, g, i, m, p\}$	$\{f, c, a, m, p\}$	$min_support = 3$
200	$\{a, b, c, f, l, m, o\}$	$\{f, c, a, b, m\}$	- **
300	$\{b, f, h, j, o\}$	$\{f, b\}$	
400	$\{b, c, k, s, p\}$	$\{c, b, p\}$	
500	$\{a, f, c, e, l, p, m, n\}$	$\{f, c, a, m, p\}$	

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



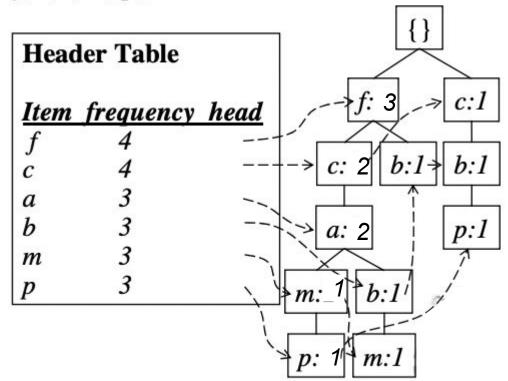
TID	Items bought	(ordered) frequent items	
100	$\{f, a, c, d, g, i, m, p\}$	$\{f, c, a, m, p\}$	$min_support = 3$
200	$\{a, b, c, f, l, m, o\}$	$\{f, c, a, b, m\}$	- 11
300	$\{b, f, h, j, o\}$	$\{f, b\}$	
400	$\{b, c, k, s, p\}$	$\{c, b, p\}$	
500	$\{a, f, c, e, l, p, m, n\}$	$\{f, c, a, m, p\}$	

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



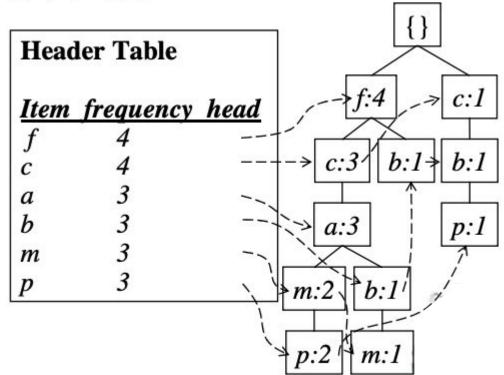
TID	Items bought	(ordered) frequent items	
100	$\{f, a, c, d, g, i, m, p\}$	$\{f, c, a, m, p\}$	$min_support = 3$
200	$\{a, b, c, f, l, m, o\}$	$\{f, c, a, b, m\}$	_ 11
300	$\{b, f, h, j, o\}$	$\{f, b\}$	
400	$\{b, c, k, s, p\}$	$\{c, b, p\}$	
500	$\{a, f, c, e, l, p, m, n\}$	$\{f, c, a, m, p\}$	

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



<b>TID</b>	Items bought	(ordered) frequent items	
100	$\{f, a, c, d, g, i, m, p\}$	$\{f, c, a, m, p\}$	$min_support = 3$
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	$\{a, f, c, e, l, p, m, n\}$	$\{f, c, a, m, p\}$	

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



## **Benefits of the FP-tree Structure**

#### Completeness:

- never breaks a long pattern of any transaction
- preserves complete information for frequent pattern mining

#### Compactness

- reduce irrelevant information—infrequent items are gone
- frequency descending ordering: more frequent items are more likely to be shared
- never be larger than the original database (if not count node-links and counts)

## Mining Frequent Patterns Using FP-tree

- General idea (divide-and-conquer)
  - Recursively grow frequent pattern path using the FP-tree

#### Method

- For each 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)

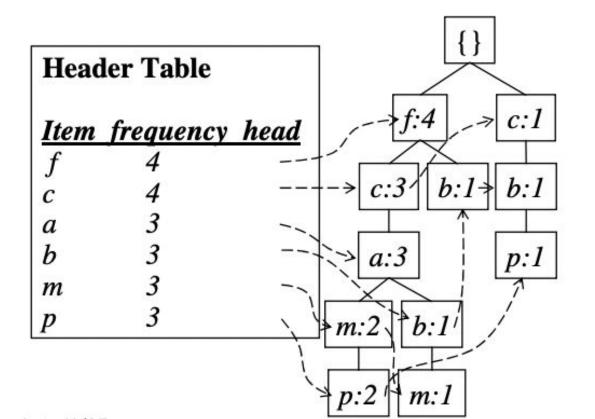
## **Major Steps to Mine FP-tree**

- 1. Construct conditional pattern base for each item in the FP-tree
- 2. Construct conditional FP-tree from each conditional pattern-base
- Recursively mine conditional FP-trees and grow frequent patterns obtained so far
- 4. If the conditional FP-tree contains a single path, simply enumerate all the patterns

#### Properties of FP-tree for Conditional Pattern Base Construction

- 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  $a_i$ .

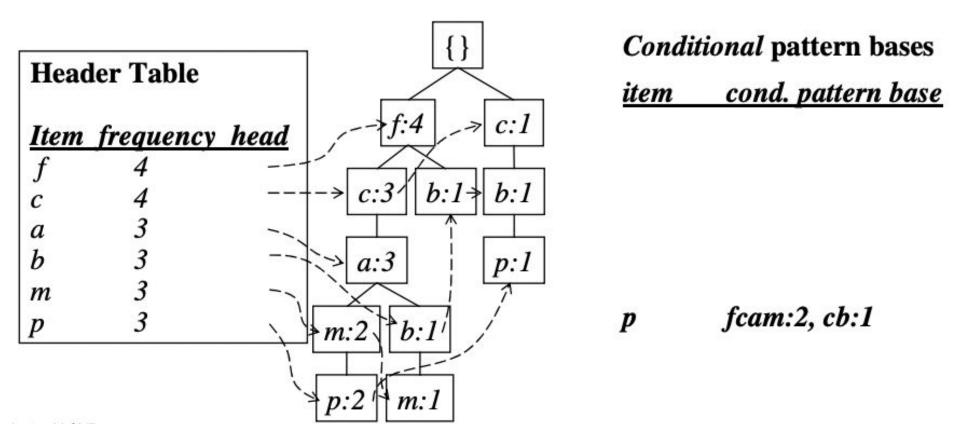
- Starting at the frequent header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item
- Accumulate all of transformed prefix paths of that item to form a conditional pattern base



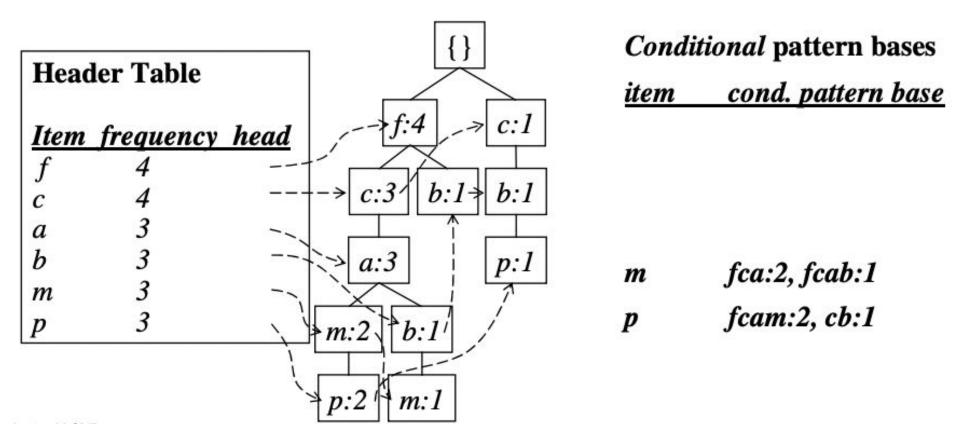
Conditional pattern bases

item cond. pattern base

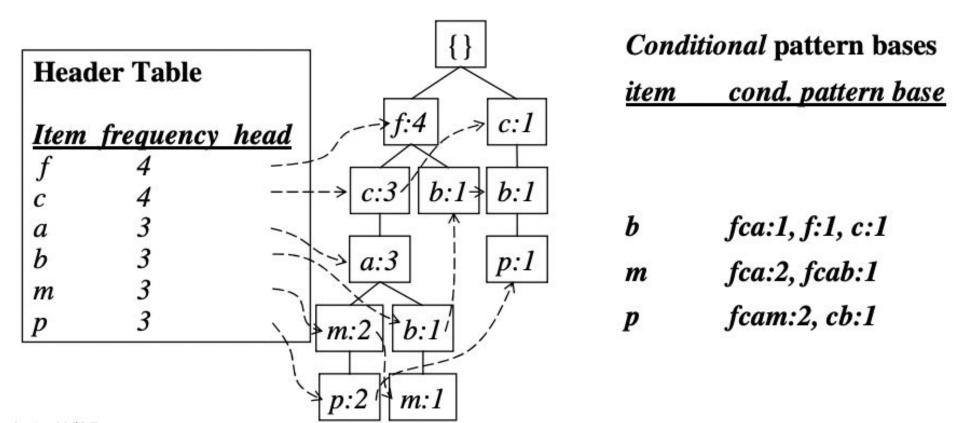
- Starting at the frequent header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item
- Accumulate all of transformed prefix paths of that item to form a conditional pattern base



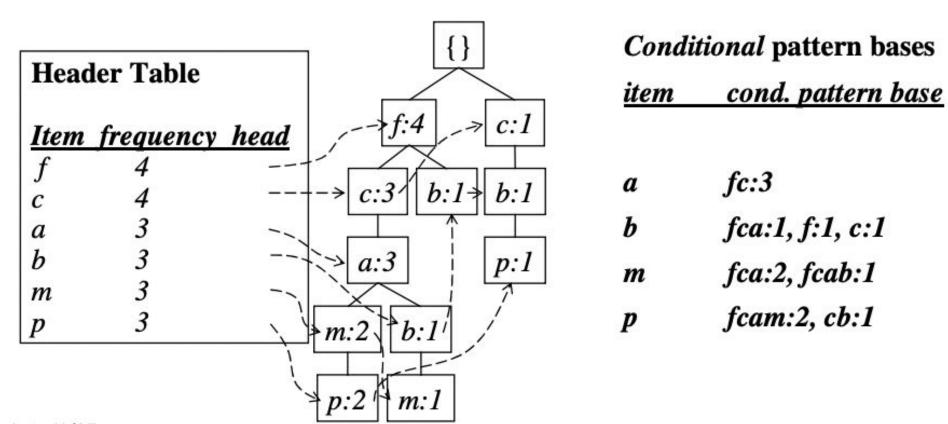
- Starting at the frequent header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item
- Accumulate all of transformed prefix paths of that item to form a conditional pattern base



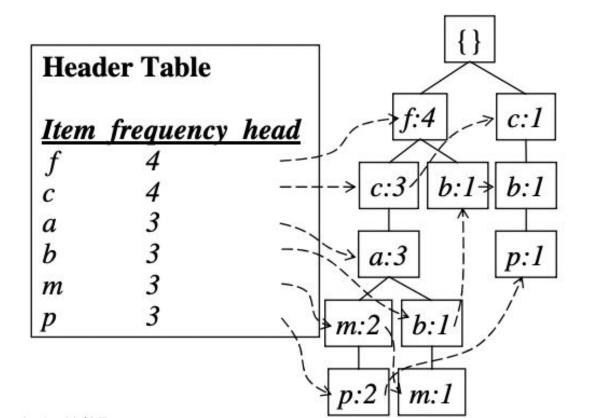
- Starting at the frequent header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item
- Accumulate all of transformed prefix paths of that item to form a conditional pattern base



- Starting at the frequent 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



- Starting at the frequent 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

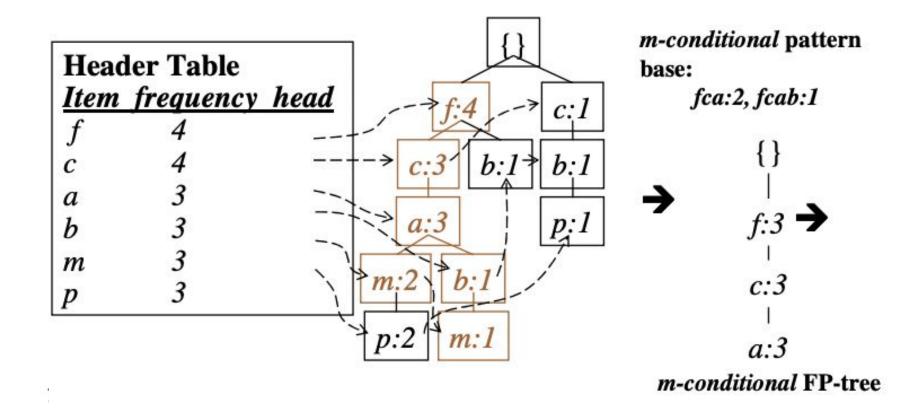


#### Conditional pattern bases

<u>item</u>	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		

# **Step 2: Construct Conditional FP-tree**

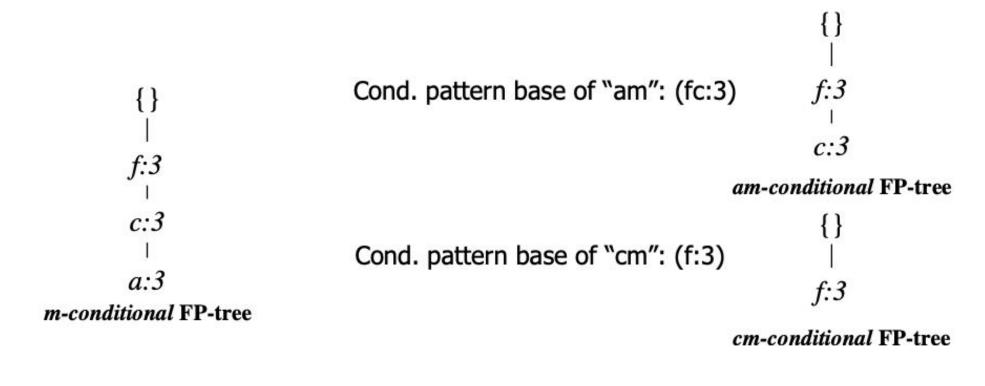
- For each pattern-base
  - Accumulate the count for each item in the base
  - Construct the FP-tree for the frequent items of the pattern base



#### Mining Frequent Patterns by Creating Conditional Pattern Bases

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
Ь	{(fca:1), (f:1), (c:1)}	Empty
α	{(fc:3)}	{(f:3, c:3)} a
С	{(f:3)}	{(f:3)} c
f	Empty	Empty

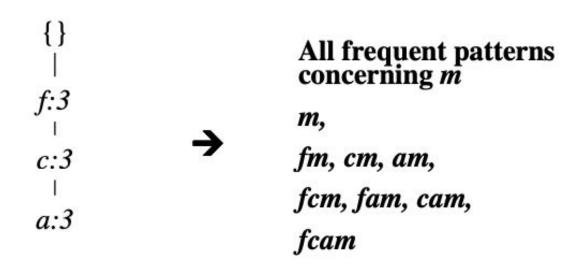
## Step 3: recursively mine the conditional FP-tree



Cond. pattern base of "cam": (f:3) f:3 cam-conditional FP-tree

# Single FP-tree Path Generation

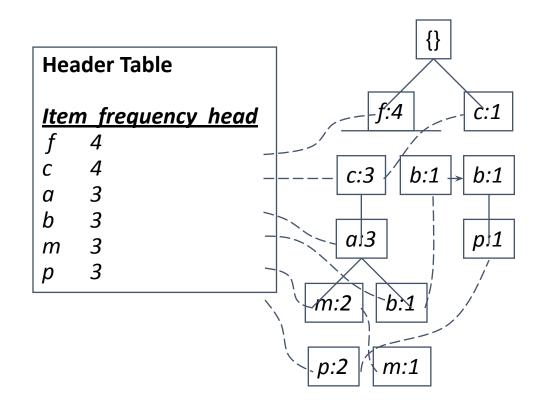
- Suppose an FP-tree T has a single path P
- The complete set of frequent pattern of *T* can be generated by enumeration of all the combinations of the sub-paths of *P*



*m-conditional* FP-tree

#### 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

#### <u>item cond. pattern base</u>

c f:3

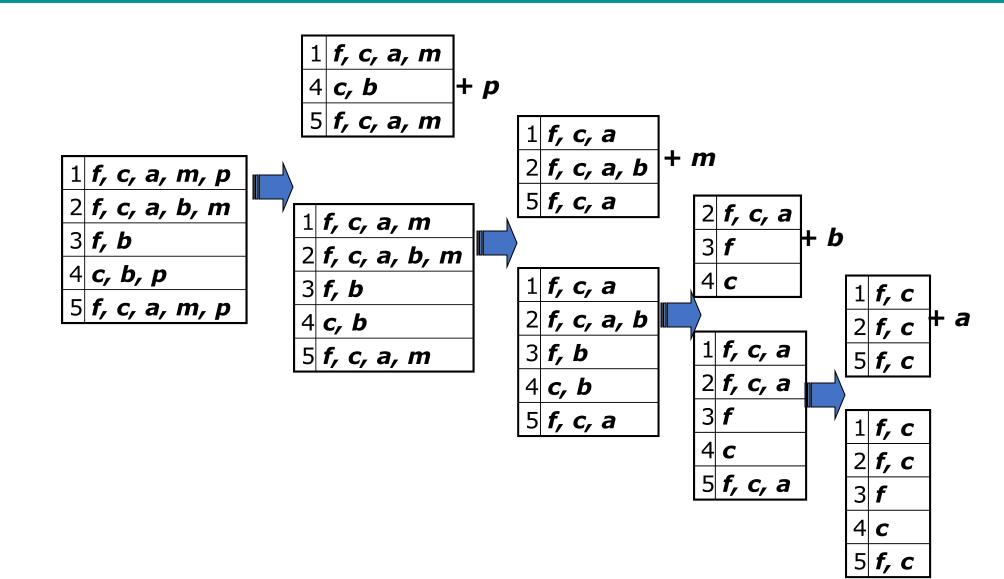
a fc:3

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

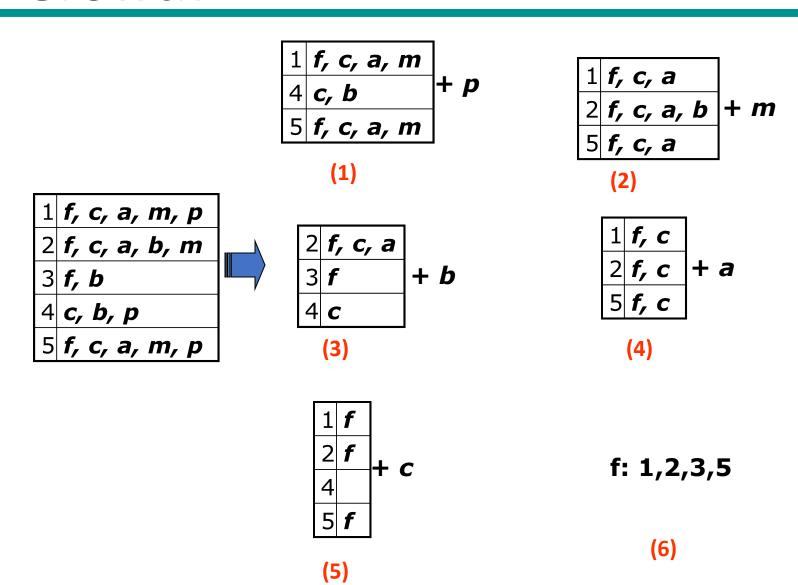
m fca:2, fcab:1

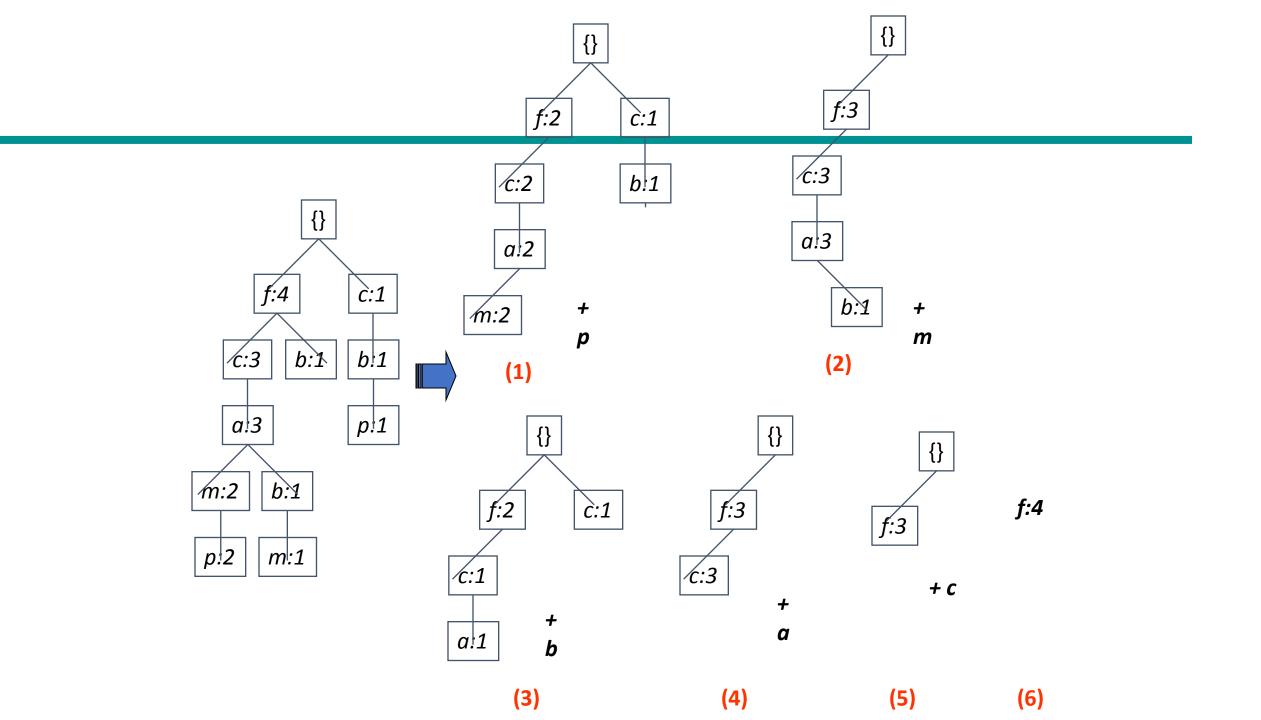
p fcam:2, cb:1

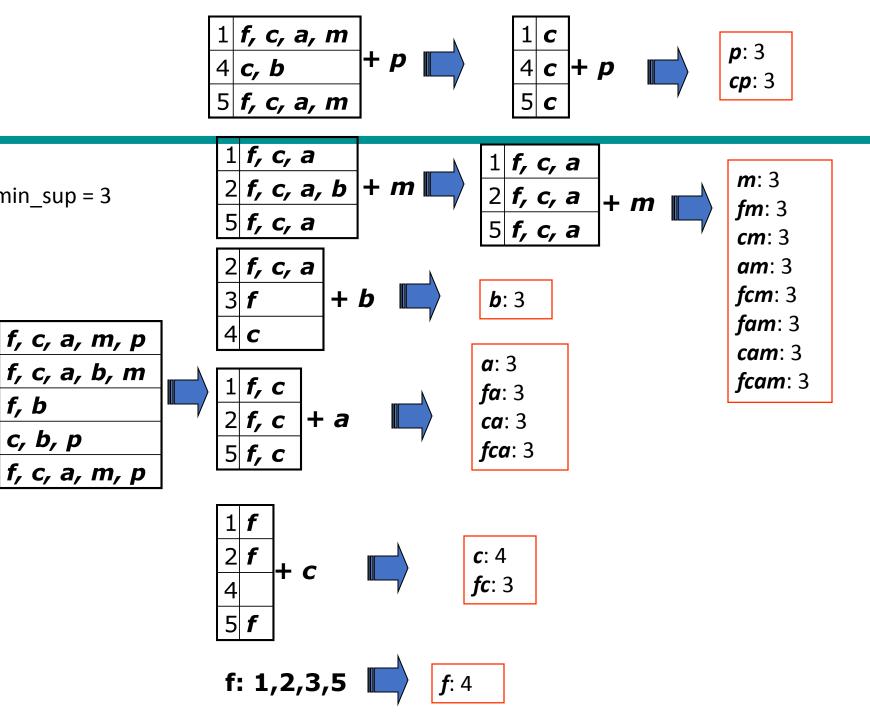
#### **FP-Growth**



#### **FP-Growth**







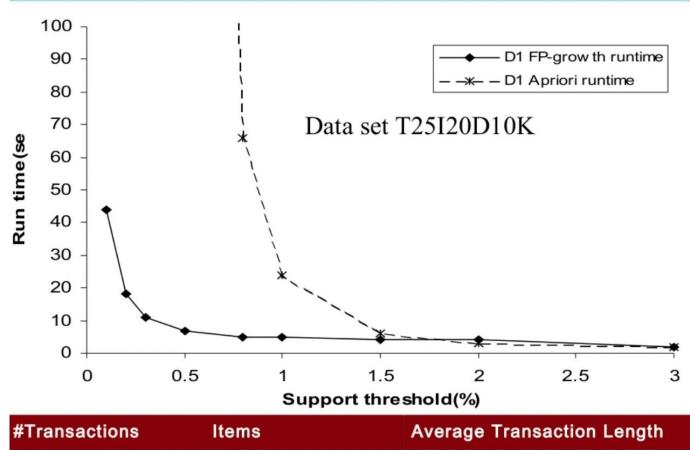
 $min_sup = 3$ 

**f, b** 

*c, b, p* 

# Why is FP-Growth Fast?

- FP-Growth is an order of magnitude faster than Apriori
  - No candidate generation, no candidate test
  - Use compact data structure
  - Eliminate repeated dataset scan
  - Basic operation is counting and FP-tree building



#Transactions	Items	Average Transaction Length
250,000	1000	12

## References

• Pattern Mining. Chapter 5. Introduction to Data Mining.

