

# Association Rules Mining



# Reference

- R. Agrawal, T. Imielinski, and A. Swami. **Mining association rules between sets of items in large databases.** In Proceedings of the ACM SIGMOD International Conference on Management of Data, pages 207-216, Washington D.C., May 1993
- R. Agrawal and R. Srikant. ***Fast algorithms for mining association rules in large databases.*** In Jorge B. Bocca, Matthias Jarke, and Carlo Zaniolo, editors, Proceedings of the 20th International Conference on Very Large Data Bases, VLDB, pages 487-499, Santiago, Chile, September 1994

# Table of Contents

---

- ▣ Introduction
- ▣ Formal statement & Decomposition
- ▣ Apriori Algorithm
- ▣ Conclusion

# Introduction

- Retail organizations have a massive amounts of sales data, referred to as the basket data.
  - From that data, how can we know the associations between items bought in a transaction?
- In 1993, Rakes Agrawal proposed a method called: Association Rules Mining.

# Introduction

Transaction ID	Items bought
1	milk, butter, juice, coke, bread
2	bread, rice, butter, milk, pasta
3	potatoes, butter, milk, bread
...	...
1988	Pepsi, bread, milk, butter, sausage

■ An association rule is:

[Butter, Bread]  [Milk]

# Formal statement & Decomposition

- ▣ Set of items:  $I = \{i_1, i_2, i_3, \dots, i_n\}$
- ▣ Transaction:  $(TransacId, T) : T \subset I$
- ▣ Itemset, k-itemset:  $X, Y \subset I$
- ▣ Association rule:  $X \Rightarrow Y : X \cap Y = \Phi$
- ▣ Support
  - ▣ Support of the itemset X: Number of transactions that contain X
  - ▣ Support of the rule  $X \Rightarrow Y$  : Support of  $X \cup Y$

# Formal statement & Decomposition

- ▣ Confidence

$$Conf(X \Rightarrow Y) = \frac{Support(X \Rightarrow Y)}{Support(X)}$$

- ▣ Association rules mining: Generate all association rules that have support and confidence greater than the user-specified minimum support (minSup) and minimum confidence (minConf) respectively.

# Formal statement & Decomposition

- Association Rules Mining can be decomposed into 2 sub problems:
- Sub problem 1: Finding all frequent itemsets.
  - $X$  is a frequent itemset  $\Leftrightarrow \text{Support}(X) \geq \text{minSup}$
  - $X$  is a frequent itemset  $\Leftrightarrow$  All subsets of  $X$  are frequent itemsets.



# Formal statement & Decomposition


- Sub problem 2: Use the frequent itemsets to generate the desired rules.

$Z$  is a frequent itemset

$$X \subseteq Z, Y \subseteq Z$$

$$X \cap Y = \Phi$$

$$Conf(X \Rightarrow Y) \geq MinConf$$

  $X \Rightarrow Y$  is an association rule

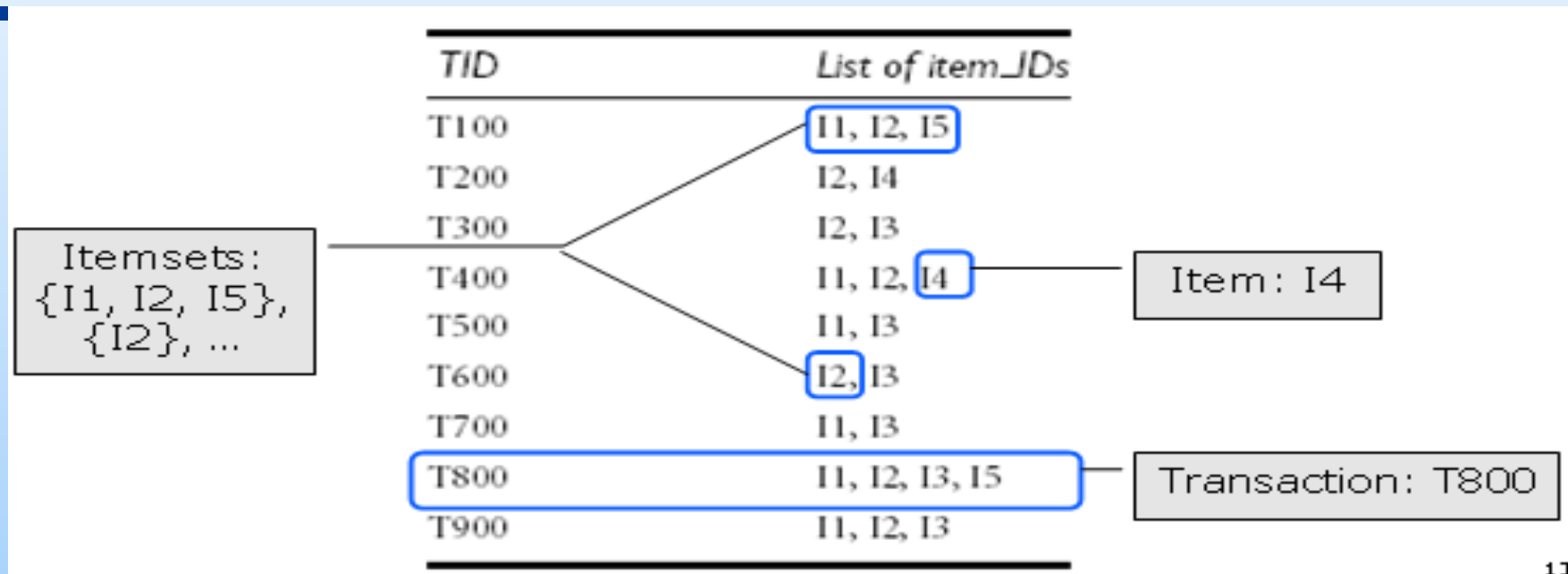
- $X$  may be find by exhaustive search.

# Outline

---

- ▣ Introduction
- ▣ Formal statement & Decomposition
- ▣ Apriori Algorithm
- ▣ Conclusion

# Apriori Algorithm - Concepts



- $X$  is **frequent** iff  $\text{support}(X) \geq \text{min\_sup}$   
Ex:  $\text{min\_sup} = 3/9$ ,  $\text{support}\{I1, I2\} = 4/9$ . So  $\{I1, I2\}$  is frequent
- $X$  is a **k-itemset** iff  $X$  has  $k$  items  
Ex:  $\{I1, I2, I5\}$  is a 3-itemset

# Apriori Algorithm - Purpose

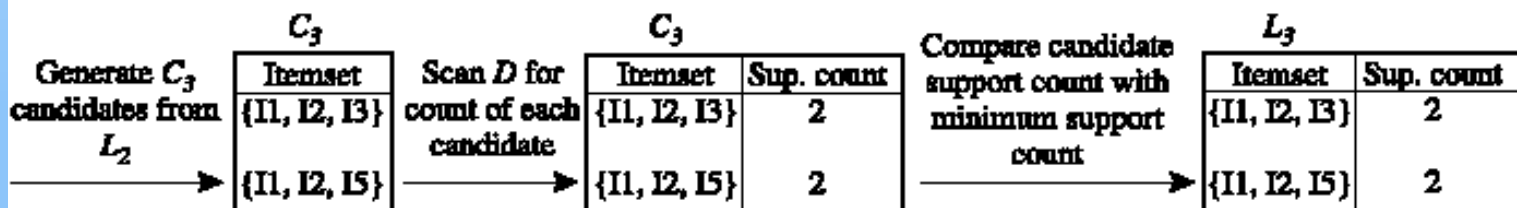
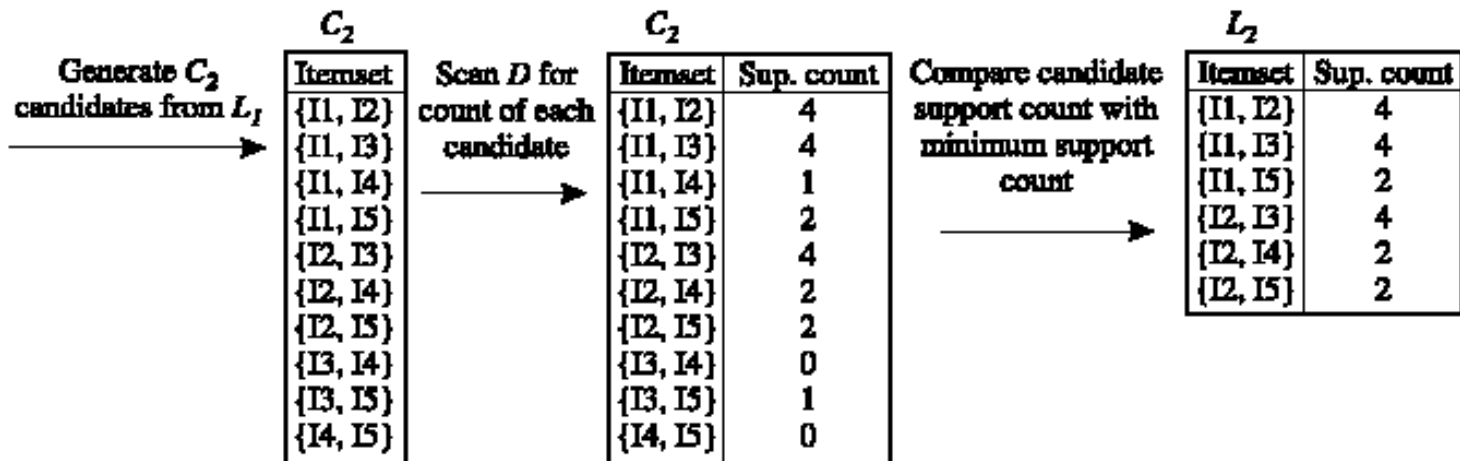
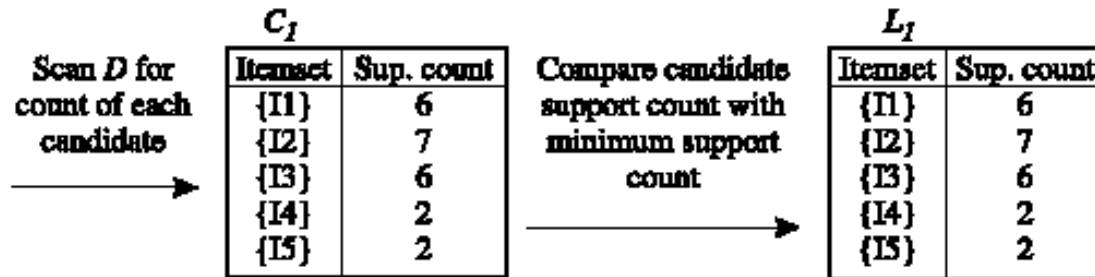
<i>TID</i>	<i>List of item_IDs</i>
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

- Find all k-itemsets having support  $\geq \text{min\_sup}$
- $k = 1, 2, \dots, 5$

# Apriori Algorithm – Step by Step

min\_sup = 2/9

minimum support count = 2



# Apriori Algorithm - Properties

- Use **prior** knowledge
- Apriori property: **All nonempty subsets of a frequent itemset must also be frequent**
- Iterative search for frequent itemsets
- Compute  **$k+1$** -itemsets from  **$k$** -itemsets

# Apriori Algorithm - Methods

Input:

- $D$ , a database of transactions;
- $min\_sup$ , the minimum support count threshold.

Output:  $L$ , frequent itemsets in  $D$ .

Method:

```
(1)   $L_1 = \text{find\_frequent\_1-itemsets}(D)$ ;  
(2)  for ( $k = 2; L_{k-1} \neq \phi; k++$ ) {  
(3)     $C_k = \text{apriori\_gen}(L_{k-1})$ ;  
(4)    for each transaction  $t \in D$  { // scan  $D$  for counts  
(5)       $C_t = \text{subset}(C_k, t)$ ; // get the subsets of  $t$  that are candidates  
(6)      for each candidate  $c \in C_t$   
(7)         $c.\text{count}++$ ;  
(8)    }  
(9)     $L_k = \{c \in C_k \mid c.\text{count} \geq min\_sup\}$   
(10) }  
(11) return  $L = \cup_k L_k$ ;
```

# Apriori Algorithm - Methods

procedure apriori\_gen( $L_{k-1}$ :frequent  $(k-1)$ -itemsets)

```
(1)   for each itemset  $l_1 \in L_{k-1}$ 
(2)     for each itemset  $l_2 \in L_{k-1}$ 
(3)       if  $(l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2]) \wedge \dots \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$  then {
(4)          $c = l_1 \bowtie l_2$ ; // join step: generate candidates
(5)         if has_infrequent_subset( $c, L_{k-1}$ ) then
(6)           delete  $c$ ; // prune step: remove unfruitful candidate
(7)         else add  $c$  to  $C_k$ ;
(8)       }
(9)   return  $C_k$ ;
```

procedure has\_infrequent\_subset( $c$ : candidate  $k$ -itemset;

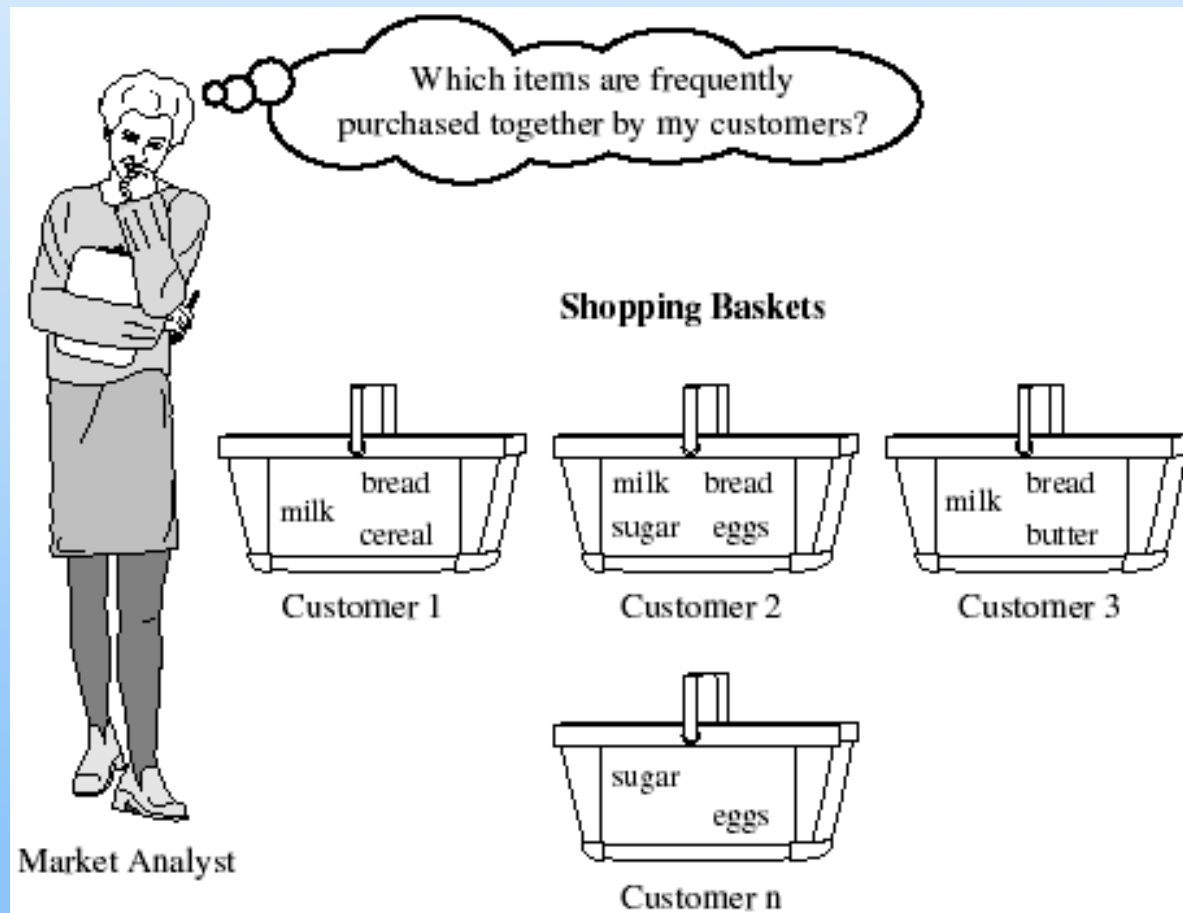
$L_{k-1}$ : frequent  $(k-1)$ -itemsets); // use prior knowledge

```
(1)   for each  $(k-1)$ -subset  $s$  of  $c$ 
(2)     if  $s \notin L_{k-1}$  then
(3)       return TRUE;
(4)   return FALSE;
```



# Conclusion

## ■ Case Study 1 – Market Basket Analysis



# Case Study 2 – Related Products

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 Links

## Customers Who Bought This Item Also Bought

		
<a href="#">Management Information Systems: Managing the...</a> by Kenneth C. Laudon	<a href="#">The Internet Book: Everything You Need to Kno...</a> by Douglas E. Comer	<a href="#">Project Management with MS Project CD + Student...</a> by Erik W. Larson
★★★★☆ (5)	★★★★☆ (10) \$47.53	★★★★☆ (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

# Other applications

---

- ▣ Catalog design
- ▣ Classification
- ▣ Clustering