

# Chapter 5

## Frequent Patterns and Association Rule Mining

1

### Outline

- Frequent Itemsets and Association Rule
- APRIORI
- Post-processing
- Applications

2

## Transactional Data

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

- Definitions:
  - An *item*: an article in a basket, or an attribute-value pair
  - A *transaction*: set of items purchased in a basket; it may have TID (transaction ID)
  - A *transactional dataset*: A set of transactions

3

## Itemsets & Frequent Itemsets

- **itemset**: A set of one or more items
- **k-itemset**  $X = \{x_1, \dots, x_k\}$
- **(absolute) support**, or, **support count** of  $X$ : Frequency or occurrence of an itemset  $X$
- **(relative) support**,  $s$ , is the fraction of transactions that contains  $X$  (i.e., the **probability** that a transaction contains  $X$ )
- An itemset  $X$  is **frequent** if  $X$ 's support is no less than a **minsup** threshold

4

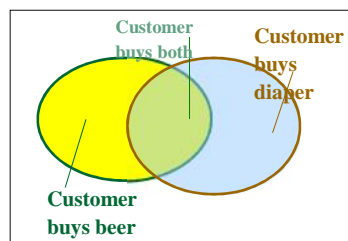
## Association Rule

- Find all the rules  $X \rightarrow Y$  with minimum support and confidence
  - **support**,  $s$ , **probability** that a transaction contains  $X \cup Y$ 
    - $s = P(X \hat{=} Y)$   
 $= \text{support count } (X \hat{=} Y) / \text{number of all transactions}$
  - **confidence**,  $c$ , **conditional probability** that a transaction having  $X$  also contains  $Y$ 
    - $c = P(X|Y)$   
 $= \text{support count } (X \cup Y) / \text{support count } (X)$

5

## Example

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk



- Let  $\text{minsup} = 50\%$ ,  $\text{minconf} = 50\%$ 
  - Number of all transactions = 5 & Min. support count =  $5 \times 50\% = 2.5 \rightarrow 3$
  - Items: Beer, Nuts, Diaper, Coffee, Eggs, Milk
  - Freq. Pat.: {Beer}:3, {Nuts}:3, {Diaper}:4, {Eggs}:3, {Beer, Diaper}:3
- Association rules (support, confidence):
  - $\text{Beer} \rightarrow \text{Diaper}$  (60%, 100%)
  - $\text{Diaper} \rightarrow \text{Beer}$  (60%, 75%)

6

## Use of Association Rules

- Association rules do not necessarily represent causality or correlation between the two itemsets.
  - $X \not\Rightarrow Y$  does not mean  $X$  causes  $Y$ , no Causality
  - $X \not\Rightarrow Y$  can be different from  $Y \not\Rightarrow X$ , unlike correlation
- Association rules assist in Basket data analysis, cross-marketing, catalog design, sale campaign analysis, Web log (click stream) analysis, and DNA sequence analysis.
  - What products were often purchased together?— Beer and diapers?!
  - What are the subsequent purchases after buying a PC?
  - What kinds of DNA are sensitive to this new drug?
  - Can we automatically classify web documents?

7

## Computational Complexity of Frequent Itemset Mining

- How many itemsets are potentially to be generated in the worst case?
  - The number of frequent itemsets to be generated is sensitive to the minsup threshold
  - When minsup is low, there exist potentially an exponential number of frequent itemsets
  - The worst case is close to  $M^N$  where  $M$ : # distinct items, and  $N$ : max length of transactions when  $M$  is large.
- The worst case complexity vs. the expected probability
  - Ex. Suppose Walmart has  $10^4$  kinds of products
    - The chance to pick up one product  $10^{-4}$
    - The chance to pick up a particular set of 10 products:  $\sim 10^{-40}$
    - What is the chance this particular set of 10 products to be frequent  $10^3$  times in  $10^9$  transactions?

8

## Outline

- Frequent Itemsets and Association Rule
- **APRIORI**
- Post-processing
- Applications

9

## Association Rule Mining

- Major steps in association rule mining
  - Frequent itemsets computation
  - Rule derivation
- Use of support and confidence to measure strength

10

## APRIORI

### ■ Method:

- Initially, scan DB once to get frequent 1-itemset
- **Generate** length  $(k+1)$  **candidate** itemsets ( $C_{k+1}$ ) from length  $k$  **frequent** itemsets ( $L_k$ )
- **Test** the candidates against DB
- Terminate when no frequent or candidate set can be generated

11

## Implementation of Apriori

- ### ■ How to generate candidates $C_k$ ?
- Step 1: self-joining  $L_{k-1}$
  - Step 2: pruning

12

## Self-joining

- $C_k$  is generated by joining  $L_{k-1}$  with itself
  - Apriori assumes that items within a transaction or itemset are sorted in lexicographic order.
  - $l_1$  from  $L_{k-1} = \{l_1[1], l_1[2], \dots, l_1[k-2], l_1[k-1]\}$
  - $l_2$  from  $L_{k-1} = \{l_2[1], l_2[2], \dots, l_2[k-2], l_2[k-1]\}$
  - Only when first **k-2** items of  $l_1$  and  $l_2$  are in common, and  $l_1[k-1] < l_2[k-1]$ , these two itemsets are joinable

13

## Example of Self-joining

- $L_3 = \{\{a,b,c\}, \{a,b,d\}, \{a,c,d\}, \{a,c,e\}, \{b,c,d\}\}$
- Self-joining:  $L_3 * L_3$ 
  - $abcd$  from  $abc$  and  $abd$
  - $acde$  from  $acd$  and  $ace$

14

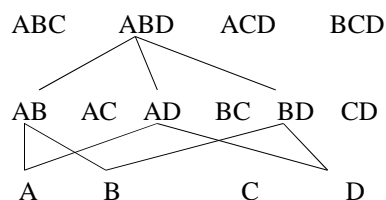
## The Downward Closure Property

- The **downward closure** property of frequent patterns
  - Any subset of a frequent itemset must be frequent
  - If **{beer, diaper, nuts}** is frequent, so is **{beer, diaper}**
  - i.e., every transaction having {beer, diaper, nuts} also contains {beer, diaper}

15

## Apriori: A Candidate Generation & Test Approach

- A *frequent* (used to be called large) *itemset* is an itemset whose support (S) is  $\geq \text{minSup}$ .
- **Apriori pruning principle:**
  - If there is **any** itemset which is infrequent, its superset should not be generated/tested!



16

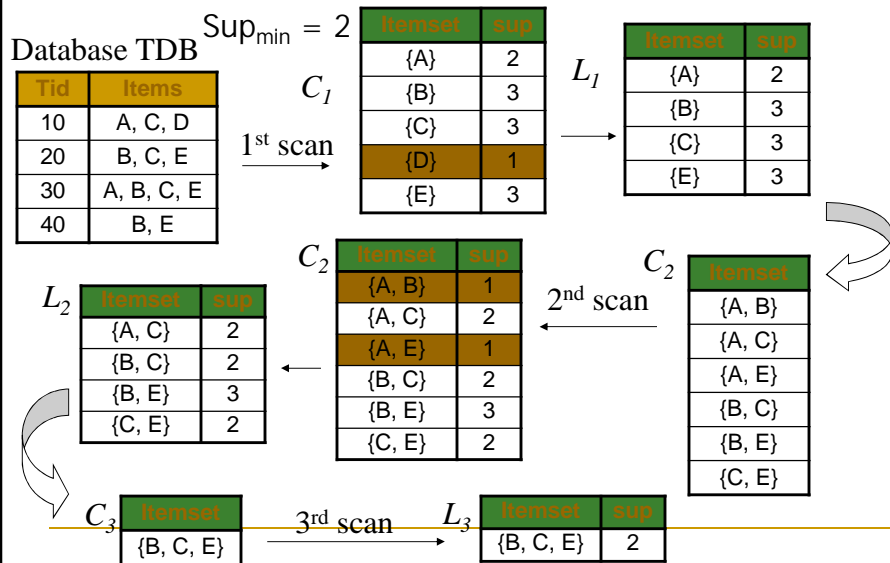


## Pruning Example

- $L_3 = \{\{a,b,c\}, \{a,b,d\}, \{a,c,d\}, \{a,c,e\}, \{b,c,d\}\}$
- Self-joining:  $L_3 * L_3$ 
  - $abcd$  from  $abc$  and  $abd$
  - $acde$  from  $acd$  and  $ace$
- Pruning:
  - $abcd$  is kept since  $abc$ ,  $abd$ ,  $acd$ , and  $bcd$  are in  $L_3$
  - $acde$  is removed because  $cde$  and  $ade$  are not in  $L_3$
- $C_4 = \{abcd\}$

17

## The Apriori Algorithm—An Example



18

## The Apriori Algorithm (Pseudo-Code)

$C_k$ : Candidate itemset of size  $k$

$L_k$ : frequent itemset of size  $k$

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

**for** ( $k = 1; L_k \neq \emptyset; k++$ ) **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**  $\cup_k L_k$ ;

19

## Further Improvement of the Apriori Method

- Major computational challenges
  - Multiple scans of transaction database
  - Huge number of candidates
  - Tedious workload of support counting for candidates
- Improving Apriori: general ideas
  - Reduce passes of transaction database scans
  - Shrink number of candidates
  - Facilitate support counting of candidates
- Completeness: any association rule mining algorithm should get the same set of frequent itemsets.

20

## Partition: Scan Database Only Twice

- Any itemset that is potentially frequent in DB must be frequent in at least one of the partitions of DB
  - Scan 1: partition database and find local frequent patterns
    - Since the sub-database is relatively smaller than original database, all frequent itemsets are tested in one scan.
  - Scan 2: consolidate global frequent patterns
- A. Savasere, E. Omiecinski, and S. Navathe. *An efficient algorithm for mining association in large databases*. In *VLDB'95*

21

## DHP: Reduce the Number of Candidates

- A  $k$ -itemset whose corresponding hashing bucket count is below the threshold cannot be frequent
  - Candidates: a, b, c, d, e
  - Frequent 1-itemset: a, b, d, e
  - Using some hash function, get hash entries: {ab, ad, ae} with bucket count as 3, {bd, bd, be, de} with bucket count as 4...
  - ab is not a candidate 2-itemset if the count of bucket - {ab, ad, ae} is below support threshold
- J. Park, M. Chen, and P. Yu. *An effective hash-based algorithm for mining association rules*. In *SIGMOD'95*

22

## Transaction Reduction

- A transaction that does not contain any frequent  $k$ -itemsets cannot contain any frequent  $(k+1)$ -itemsets.
- Such a transaction could be marked or removed from further consideration of  $(k+1)$ -itemsets.
  - Less support counting

23

## Sampling for Frequent Patterns

- Select a sample of original database, mine frequent patterns within sample using Apriori with a lower support count threshold than min. support.
- Scan database once to verify frequent itemsets found in sample, only *borders* of closure of frequent patterns are checked
  - Example: check *abcd* instead of *ab*, *ac*, ..., etc.
- Scan database again to find missed frequent patterns
- H. Toivonen. Sampling large databases for association rules. In *VLDB'96*

24

## Outline

- Frequent Itemsets and Association Rule
- APRIORI
- Post-processing
- Applications

## Derive rules from frequent itemsets

- Frequent itemsets  $\neq$  association rules
- One more step is required to find association rules
- For each frequent itemset  $X$ ,  
For each proper nonempty subset  $A$  of  $X$ ,
  - Let  $B = X - A$
  - $A \Rightarrow B$  is an association rule if
    - Confidence  $(A \Rightarrow B) \geq \text{minConf}$ ,  
where  $\text{support}(A \Rightarrow B) = \text{support}(AB)$  and  
confidence  $(A \Rightarrow B) = \text{support}(AB) / \text{support}(A)$

### Example – deriving rules from frequent itemsets

- Suppose {2,3,4} is frequent, with supp=50%
  - Proper nonempty subsets: {2,3}, {2,4}, {3,4}, {2}, {3}, {4}, with supp = 50%, 50%, 75%, 75%, 75%, 75% respectively
  - These generate these association rules:
    - 2,3  $\Rightarrow$  4, confidence=100%
    - 2,4  $\Rightarrow$  3, confidence=100%
    - 3,4  $\Rightarrow$  2, confidence=67%
    - 2  $\Rightarrow$  3,4, confidence=67%
    - 3  $\Rightarrow$  2,4, confidence=67%
    - 4  $\Rightarrow$  2,3, confidence=67%
- All rules have support = 50%

### Mining Various Kinds of Association Rules

- Mining multilevel association
- Mining multidimensional association
- Mining quantitative association

## Mining Multiple-Level Association Rules

- Items often form hierarchies
- A Top-down strategy is employed, any algorithm could be used for mining at each level.
- Flexible support settings
  - Uniform minimum support for all levels.
  - Items at the lower level are expected to have lower support.
  - Set up user-specific, item or group-based minimum support.
    - Setting particularly low support thresholds for laptop computers and flash drives.

## Example

uniform support

Level 1  
min\_sup = 5%

Level 2  
min\_sup = 5%

Milk  
[support = 10%]

2% Milk  
[support = 6%]

Skim Milk  
[support = 4%]

reduced support

Level 1  
min\_sup = 5%

Level 2  
min\_sup = 3%

## Mining Multi-Dimensional Association

- Single-dimensional rules:  
 $\text{buys}(X, \text{"milk"}) \Rightarrow \text{buys}(X, \text{"bread"})$
- Multi-dimensional rules:  $\geq 2$  dimensions or predicates
  - Inter-dimension assoc. rules (*no repeated predicates*)  
 $\text{age}(X, \text{"19-25"}) \wedge \text{occupation}(X, \text{"student"}) \Rightarrow \text{buys}(X, \text{"coke"})$
  - hybrid-dimension assoc. rules (*repeated predicates*)  
 $\text{age}(X, \text{"19-25"}) \wedge \text{buys}(X, \text{"popcorn"}) \Rightarrow \text{buys}(X, \text{"coke"})$

## Mining Quantitative Associations

- Categorical Attributes and Quantitative Attributes
- Techniques can be categorized by how numerical attributes, such as **age** or **salary** are treated
  - Static discretization based on predefined concept hierarchies
  - Dynamic discretization based on data distribution.
  - Clustering: Distance-based association.
    - One dimensional clustering then association



Id	Age	Income	Student	Credit_Rating	Buy_Computer
1	27	75,000	No	600	No
2	25	72,000	No	730	No
3	33	88,000	No	640	Yes
4	50	55,000	No	620	Yes
5	52	34,000	Yes	640	Yes
6	45	30,000	Yes	720	No
7	32	25,000	Yes	740	Yes
8	25	54,000	No	630	No
9	22	35,000	Yes	640	Yes
10	48	67,000	Yes	660	Yes
11	24	64,000	Yes	715	Yes
12	37	62,000	No	710	Yes
13	33	90,000	Yes	650	Yes
14	45	59,000	No	705	No

age	income	student	credit_rating	com
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

age	income	student	credit_rating	buy_computer
1	3	0	1	0
1	3	0	2	0
2	3	0	1	1
3	2	0	1	1
3	1	1	1	1
3	1	1	2	0
2	1	1	2	1
1	2	0	1	0
1	1	1	1	1
3	2	1	1	1
1	2	1	2	1
2	2	0	2	1
2	3	1	1	1
3	2	0	2	0

## Interestingness Measure: Correlations (Lift)

- The occurrence of itemset A is independent of the occurrence of itemset B if  $P(A \cup B) = P(A)P(B)$ ; otherwise, itemsets A and B are dependent and correlated as events.
- Measure of dependent/correlated events: **lift**

$$lift = \frac{P(A \cup B)}{P(A)P(B)}$$

- $< 1$ , A is negatively correlated with B.
- $> 1$ , A and B are positively correlated.
- $= 1$ , A and B are independent.

## Interestingness Measure: Correlations (Lift)

- *play basketball*  $\Rightarrow$  *eat cereal* [40%, 66.7%] is misleading
  - The overall % of students eating cereal is 75% > 66.7%.
- *play basketball*  $\Rightarrow$  *not eat cereal* [20%, 33.3%] is more accurate, although with lower support and confidence

$$\text{lift}(B, C) = \frac{2000 / 5000}{3000 / 5000 * 3750 / 5000} = 0.89$$

$$\text{lift}(B, \neg C) = \frac{1000 / 5000}{3000 / 5000 * 1250 / 5000} = 1.33$$

	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

## Reference

- “The application of Association Rule Mining to Remotely Sensed Data”, J. Dong, W. Perrizo, Q. Ding and J. Zhou, SAC’2000.