

Outline

- Introduction
- Basic concepts
- Frequent Itemsets Mining (FIM) – Apriori
- Association Rules Mining

Association Rules Mining

- (Recall the) 2-step method to extract the association rules:
 - Determine the frequent itemsets w.r.t. min support s ← FIM problem (Apriori)
 - Generate the association rules w.r.t. min confidence c .

- Regarding step 2, the following method is followed:
 - For every frequent itemset X
 - for every subset Y of X : $Y \neq \emptyset$, $Y \neq X$, the rule $Y \rightarrow (X - Y)$ is formed
 - Remove rules that violate min confidence c

$$\text{confidence}(Y \rightarrow (X - Y)) = \frac{\text{support_count}(X)}{\text{support_count}(Y)}$$

- Store the frequent itemsets and their supports in a hash table
 - no database access!

Let $X = \{1, 2, 3\}$ be frequent

There are 6 candidate rules that can be generated from X :

- $\{1, 2\} \rightarrow 3$
- $\{1, 3\} \rightarrow 2$
- $\{2, 3\} \rightarrow 1$
- $\{1\} \rightarrow \{2, 3\}$
- $\{2\} \rightarrow \{1, 3\}$
- $\{3\} \rightarrow \{1, 2\}$

To identify strong rules, we can use the support counts (already computed during the FIM step)

Pseudocode

Input:

D //Database of transactions
 I //Items
 L //Large itemsets
 s //Support
 α //Confidence

Output:

R //Association Rules satisfying s and α

ARGen Algorithm:

```
 $R = \emptyset;$   
for each  $l \in L$  do  
  for each  $x \subset l$  such that  $x \neq \emptyset$  and  $x \neq l$  do  
    if  $\frac{\text{support}(l)}{\text{support}(x)} \geq \alpha$  then  
       $R = R \cup \{x \Rightarrow (l - x)\};$ 
```

Confidence-based pruning

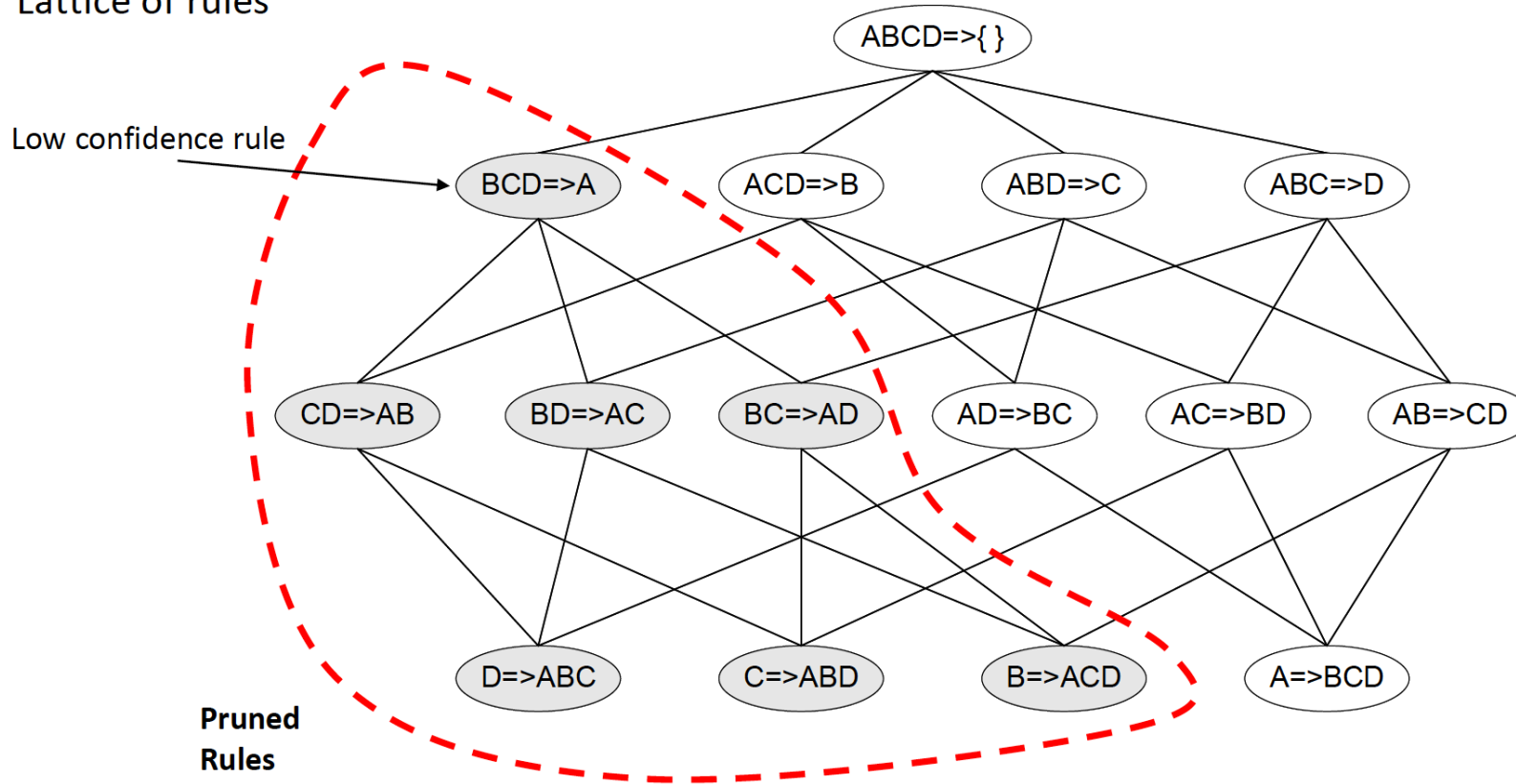
- How to efficiently generate rules from frequent itemsets?
- Confidence does not follow the monotonicity property
 - i.e., confidence $(X \rightarrow Y)$ can be $>, <, =$ to confidence $(X' \rightarrow Y')$, $X' \subseteq X, Y' \subseteq Y$
 - e.g., confidence $(ABC \rightarrow D)$ can be larger or smaller than confidence $(AB \rightarrow D)$
- But the confidence of rules generated from the same itemset does

If rule $X \rightarrow Y - X$ does not satisfy the minConfidence threshold, then any rule $X' \rightarrow Y - X'$, where $X' \subseteq X$, must not satisfy the minConfidence threshold as well.

- For example, for $X = \{ABCD\}$, then
 - confidence $(ABC \rightarrow D) \geq$ confidence $(AB \rightarrow CD) \geq$ confidence $(A \rightarrow BCD)$

Confidence-based pruning

- Lattice of rules



Example

<i>tid</i>	X_T
1	{Bier, Chips, Wine}
2	{Bier, Chips}
3	{Pizza, Wine}
4	{Chips, Pizza}

Transaction database

$I = \{\text{Bier, Chips, Pizza, Wine}\}$

Itemset	Cover	Sup.	Freq.
{}	{1,2,3,4}	4	100 %
{Bier}	{1,2}	2	50 %
{Chips}	{1,2,4}	3	75 %
{Pizza}	{3,4}	2	50 %
{Wine}	{1,3}	2	50 %
{Bier, Chips}	{1,2}	2	50 %
{Bier, Wine}	{1}	1	25 %
{Chips, Pizza}	{4}	1	25 %
{Chips, Wine}	{1}	1	25 %
{Pizza, Wine}	{3}	1	25 %
{Bier, Chips, Wine}	{1}	1	25 %

Rule	Sup.	Freq.	Conf.
$\{\text{Bier}\} \Rightarrow \{\text{Chips}\}$	2	50 %	100 %
$\{\text{Bier}\} \Rightarrow \{\text{Wine}\}$	1	25 %	50 %
$\{\text{Chips}\} \Rightarrow \{\text{Bier}\}$	2	50 %	66 %
$\{\text{Pizza}\} \Rightarrow \{\text{Chips}\}$	1	25 %	50 %
$\{\text{Pizza}\} \Rightarrow \{\text{Wine}\}$	1	25 %	50 %
$\{\text{Wine}\} \Rightarrow \{\text{Bier}\}$	1	25 %	50 %
$\{\text{Wine}\} \Rightarrow \{\text{Chips}\}$	1	25 %	50 %
$\{\text{Wine}\} \Rightarrow \{\text{Pizza}\}$	1	25 %	50 %
$\{\text{Bier, Chips}\} \Rightarrow \{\text{Wine}\}$	1	25 %	50 %
$\{\text{Bier, Wine}\} \Rightarrow \{\text{Chips}\}$	1	25 %	100 %
$\{\text{Chips, Wine}\} \Rightarrow \{\text{Bier}\}$	1	25 %	100 %
$\{\text{Bier}\} \Rightarrow \{\text{Chips, Wine}\}$	1	25 %	50 %
$\{\text{Wine}\} \Rightarrow \{\text{Bier, Chips}\}$	1	25 %	50 %

Evaluating Association Rules 1/2

Interesting and misleading association rules

Example:

- Database on the behavior of students in a school with 5.000 students
- Itemsets:
 - 60% of the students play Soccer,
 - 75% of the students eat chocolate bars
 - 40% of the students play Soccer and eat chocolate bars
- Association rules: $\{\text{"Play Soccer"}\} \rightarrow \{\text{"Eat chocolate bars"}\}$, confidence = $40\%/60\% = 67\%$
 - The rule has a high confidence, however:
 $\{\text{"Eat chocolate bars"}\}$, support = 75% , regardless of whether they play soccer.
 - Thus, knowing that one is playing soccer decreases his/her probability of eating chocolate (from $75\% \rightarrow 67\%$)
 - Therefore, the rule $\{\text{"Play Soccer"}\} \rightarrow \{\text{"Eat chocolate bars"}\}$ is misleading despite its high confidence



Evaluating Association Rules 2/2

Task: Filter out misleading rules

Let $\{A\} \rightarrow \{B\}$

- Measure of “interestingness”-score of a rule:

$$interest = \frac{support(A \cup B)}{support(A)} - support(B)$$

- the higher the value the more interesting the rule is

- Measure of dependent/correlated events:

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

- the ratio of the *observed* support to that *expected* if X and Y were independent.
- Lift > 1 means that the rule is interesting, lift < 1 means that the presence of one item has negative effect on presence of other item and vice versa.

Measuring Quality of Association Rules

For a rule $A \rightarrow B$

- Support $support(A \cup B)$ $P(E_A \cap E_B)$ $E_X :=$ Event that itemset X appears in a transaction

- e.g. $support(\text{milk, bread, butter})=20\%$, i.e. 20% of the transactions contain these

- Confidence $\frac{support(A \cup B)}{support(A)}$ $\frac{P(E_A \cap E_B)}{P(E_A)}$

- e.g. $confidence(\text{milk, bread} \rightarrow \text{butter})=50\%$, i.e. 50% of the times a customer buys milk and bread, butter is bought as well.

- Lift $\frac{support(A \cup B)}{support(A)support(B)}$ $\frac{P(E_A \cap E_B)}{P(E_A)P(E_B)}$

- e.g. $lift(\text{milk, bread} \rightarrow \text{butter})=20\%/(40\%*40\%)=1.25$. the observed support is 20%, the expected (if they were independent) is 16%.