

Data Mining

Chapter 5 Association Analysis: Basic Concepts

Introduction to Data Mining, 2nd Edition
by
Tan, Steinbach, Karpatne, Kumar

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

$\{\text{Diaper}\} \rightarrow \{\text{Beer}\},$
 $\{\text{Milk, Bread}\} \rightarrow \{\text{Eggs, Coke}\},$
 $\{\text{Beer, Bread}\} \rightarrow \{\text{Milk}\},$

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. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

- **Support**

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

- **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

02/03/2018

Introduction to Data Mining

3

Definition: Association Rule

- **Association Rule**

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

- **Rule Evaluation Metrics**

- 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:

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

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

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

02/03/2018

Introduction to Data Mining

4

Association Rule Mining Task

- Given a set of transactions T , the goal of association rule mining is to find all rules having
 - support $\geq \text{minsup}$ threshold
 - confidence $\geq \text{minconf}$ 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!**

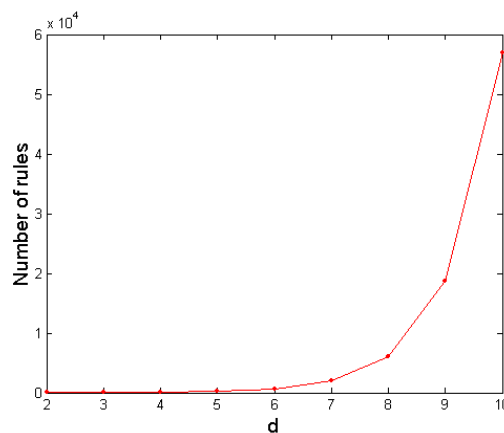
02/03/2018

Introduction to Data Mining

5

Computational Complexity

- Given d unique items:
 - Total number of itemsets = 2^d
 - Total number of possible association rules:



$$R = \sum_{k=1}^{d-1} \left[\binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right]$$
$$= 3^d - 2^{d+1} + 1$$

If $d=6$, $R = 602$ rules

02/03/2018

Introduction to Data Mining

6

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:

$\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$ ($s=0.4, c=0.67$)
 $\{\text{Milk, Beer}\} \rightarrow \{\text{Diaper}\}$ ($s=0.4, c=1.0$)
 $\{\text{Diaper, Beer}\} \rightarrow \{\text{Milk}\}$ ($s=0.4, c=0.67$)
 $\{\text{Beer}\} \rightarrow \{\text{Milk, Diaper}\}$ ($s=0.4, c=0.67$)
 $\{\text{Diaper}\} \rightarrow \{\text{Milk, Beer}\}$ ($s=0.4, c=0.5$)
 $\{\text{Milk}\} \rightarrow \{\text{Diaper, Beer}\}$ ($s=0.4, c=0.5$)

Observations:

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

02/03/2018

Introduction to Data Mining

7

Mining Association Rules

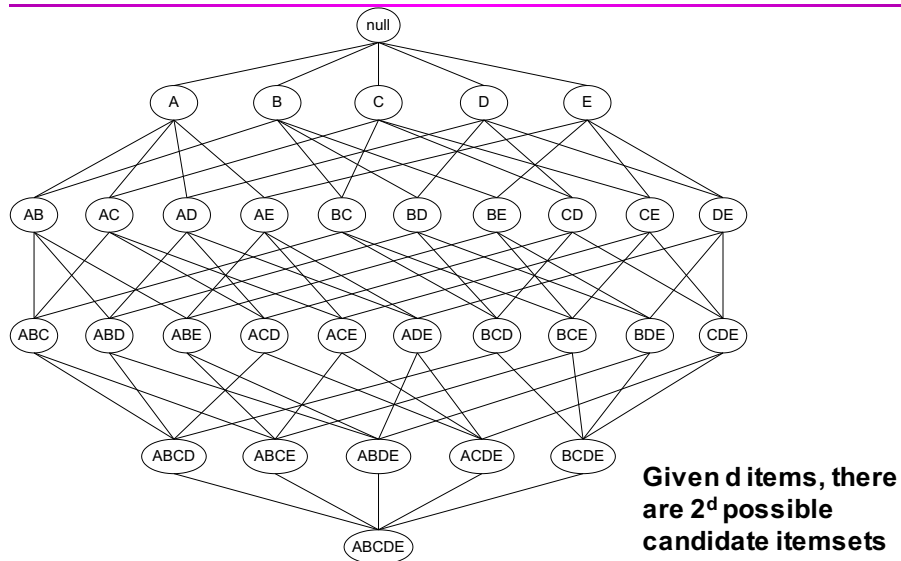
- Two-step approach:
 1. **Frequent Itemset Generation**
 - Generate all itemsets whose support \geq minsup
 2. **Rule Generation**
 - Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset
- Frequent itemset generation is still computationally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

02/03/2018

Introduction to Data Mining

8

Frequent Itemset Generation



02/03/2018

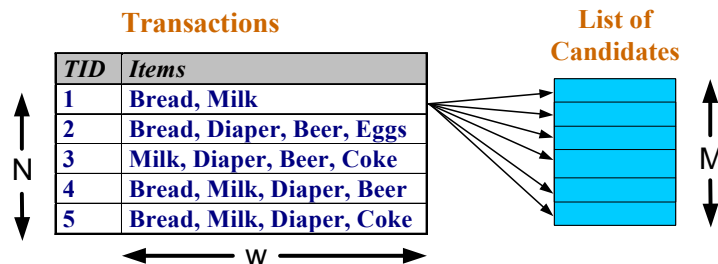
Introduction to Data Mining

9

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



- Match each transaction against every candidate
- Complexity $\sim O(NMw) \Rightarrow$ **Expensive since $M = 2^d$!!!**

02/03/2018

Introduction to Data Mining

10

Frequent Itemset Generation Strategies

- Reduce the **number of candidates** (M)
 - Complete search: $M=2^d$
 - Use pruning techniques to reduce M
- Reduce the **number of transactions** (N)
 - Reduce size of N as the size of itemset increases
 - Used by DHP and vertical-based mining algorithms
- 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

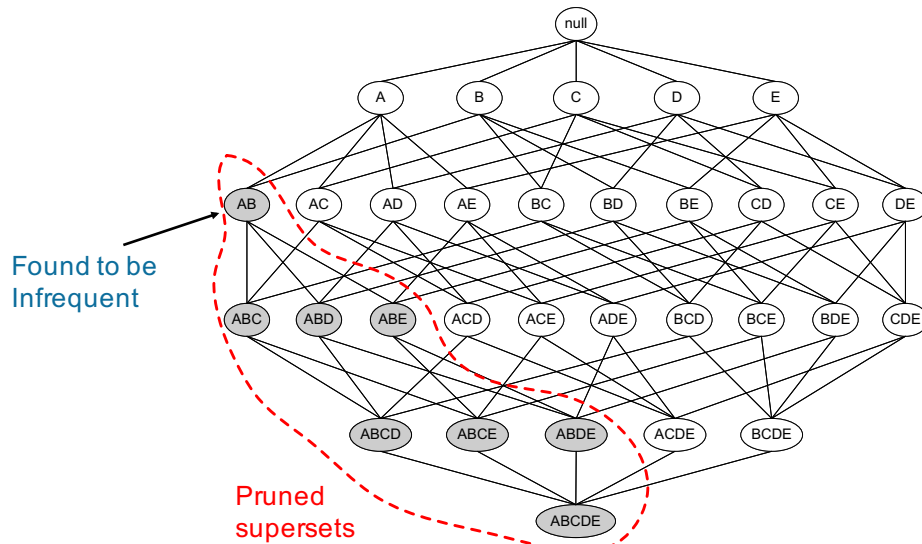
Reducing Number of Candidates

- **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:

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

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

Illustrating Apriori Principle



02/03/2018

Introduction to Data Mining

13

Illustrating Apriori Principle

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

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3$
 $6 + 15 + 20 = 41$
 With support-based pruning,
 $6 + 6 + 4 = 16$

02/03/2018

Introduction to Data Mining

14

Illustrating Apriori Principle

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

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3$
 $6 + 15 + 20 = 41$
 With support-based pruning,
 $6 + 6 + 4 = 16$

Illustrating Apriori Principle

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}

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3$
 $6 + 15 + 20 = 41$
 With support-based pruning,
 $6 + 6 + 4 = 16$

Illustrating Apriori Principle

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)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3$
 $6 + 15 + 20 = 41$
 With support-based pruning,
 $6 + 6 + 4 = 16$

Illustrating Apriori Principle

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

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3$
 $6 + 15 + 20 = 41$
 With support-based pruning,
 $6 + 6 + 4 = 16$



Triplets (3-itemsets)

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

Illustrating Apriori Principle

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

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3$
 $6 + 15 + 20 = 41$
 With support-based pruning,
 $6 + 6 + 4 = 16$
 $6 + 6 + 1 = 13$



Triplets (3-itemsets)

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

Apriori Algorithm

- F_k : frequent k-itemsets
- L_k : candidate k-itemsets

Algorithm

- Let $k=1$
- Generate $F_1 = \{\text{frequent 1-itemsets}\}$
- Repeat until F_k is empty
 - ◆ **Candidate Generation:** Generate L_{k+1} from F_k
 - ◆ **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 $\Rightarrow F_{k+1}$

Candidate Generation: Brute-force method

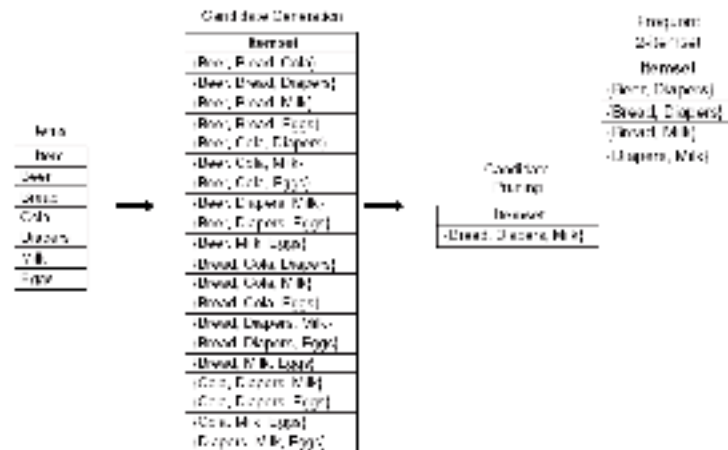


Figure 5.6. A brute-force method for generating candidate 4-itemsets.

Candidate Generation: Merge Fk-1 and F1 itemsets

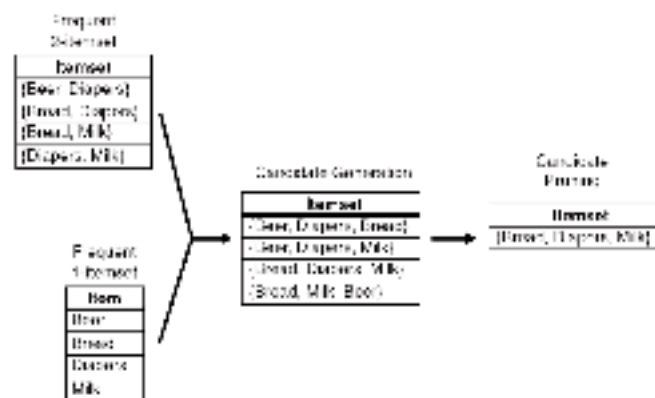


Figure 5.7. Generating and pruning candidate 4-itemsets by merging a frequent (k-1)-itemset with a frequent item. Note that some of the candidates are unnecessary because their subsets are infrequent.

Candidate Generation: Fk-1 x Fk-1 Method

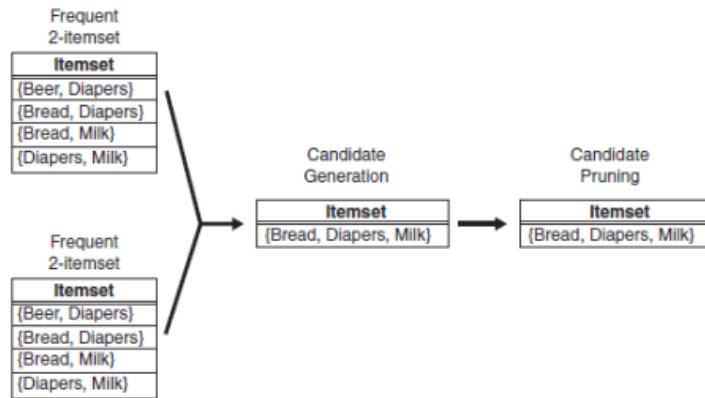


Figure 6.8. Generating and pruning candidate k -itemsets by merging pairs of frequent $(k-1)$ -itemsets.

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

- Merge two frequent $(k-1)$ -itemsets if their first $(k-2)$ items are identical
- $F_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$
 - Merge(ABC, ABD) = ABCD
 - Merge(ABC, ABE) = ABCE
 - Merge(ABD, ABE) = ABDE
 - Do not merge(ABD, ACD) because they share only prefix of length 1 instead of length 2

Candidate Pruning

- Let $F_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$ be the set of frequent 3-itemsets
- $L_4 = \{ABCD, ABCE, ABDE\}$ 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: $L_4 = \{ABCD\}$

Alternate $F_{k-1} \times F_{k-1}$ Method

- Merge two frequent (k-1)-itemsets if the last (k-2) items of the first one is identical to the first (k-2) items of the second.
- $F_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$
 - Merge(ABC, BCD) = ABCD
 - Merge(ABD, BDE) = ABDE
 - Merge(ACD, CDE) = ACDE
 - Merge(BCD, CDE) = BCDE

Candidate Pruning for Alternate $F_{k-1} \times F_{k-1}$ Method

- Let $F_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$ be the set of frequent 3-itemsets
- $L_4 = \{ABCD, ABDE, ACDE, BCDE\}$ is the set of candidate 4-itemsets generated (from previous slide)
- Candidate pruning
 - Prune ABDE because ADE is infrequent
 - Prune ACDE because ACE and ADE are infrequent
 - Prune BCDE because BCE
- After candidate pruning: $L_4 = \{ABCD\}$

02/03/2018

Introduction to Data Mining

27

Illustrating Apriori Principle

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

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)



Triplets (3-itemsets)

Itemset	Count
{Bread, Diaper, Milk}	2

Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3$
 $6 + 15 + 20 = 41$
 With support-based pruning,
 $6 + 6 + 1 = 13$

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

02/03/2018

Introduction to Data Mining

28

Support Counting of Candidate Itemsets

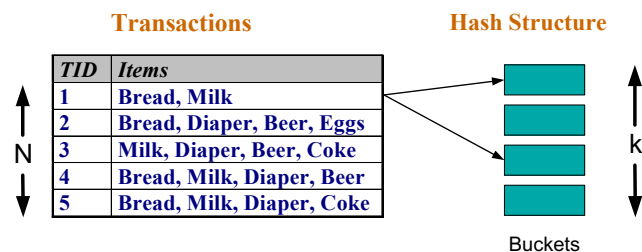
- Scan the database of transactions to determine the support of each candidate itemset
 - Must match every candidate itemset against every transaction, which is an expensive operation

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

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

Support Counting of Candidate Itemsets

- To reduce number of comparisons, store the candidate itemsets in a hash structure
 - Instead of matching each transaction against every candidate, match it against candidates contained in the hashed buckets

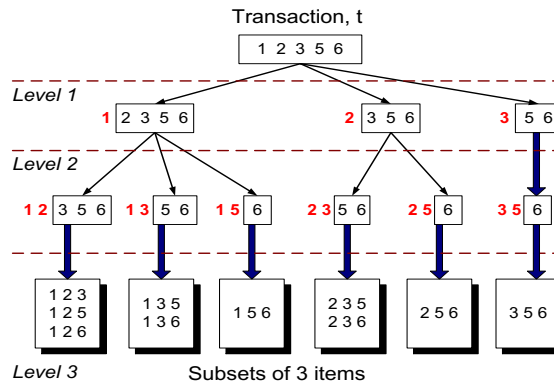


Support Counting: An Example

Suppose you have 15 candidate itemsets of length 3:

{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4}, {5 6 7}, {3 4 5},
{3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}

How many of these itemsets are supported by transaction (1,2,3,5,6)?



02/03/2018

Introduction to Data Mining

31

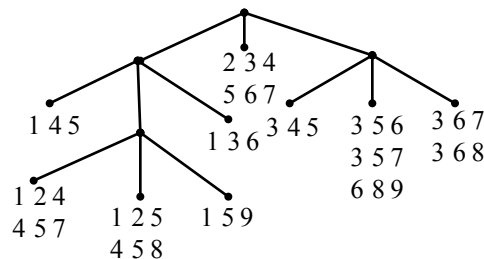
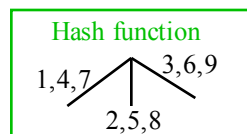
Support Counting Using a Hash Tree

Suppose you have 15 candidate itemsets of length 3:

{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4}, {5 6 7}, {3 4 5},
{3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}

You need:

- Hash function
- Max leaf size: max number of itemsets stored in a leaf node (if number of candidate itemsets exceeds max leaf size, split the node)

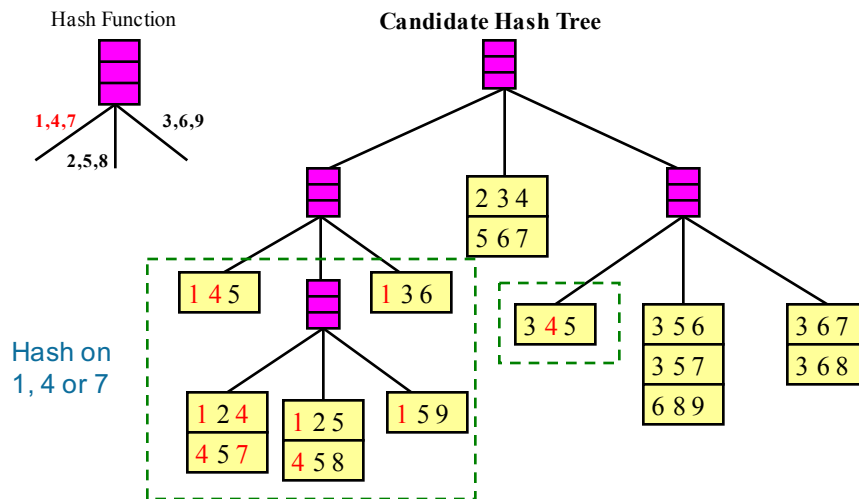


02/03/2018

Introduction to Data Mining

32

Support Counting Using a Hash Tree

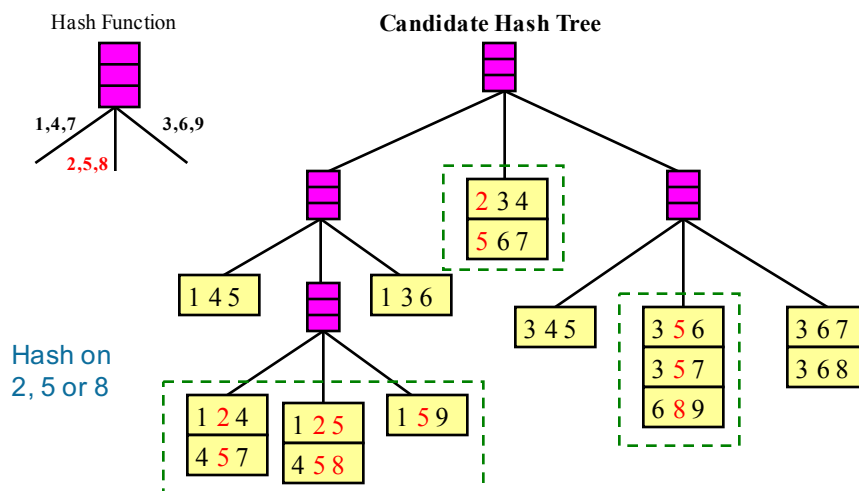


02/03/2018

Introduction to Data Mining

33

Support Counting Using a Hash Tree

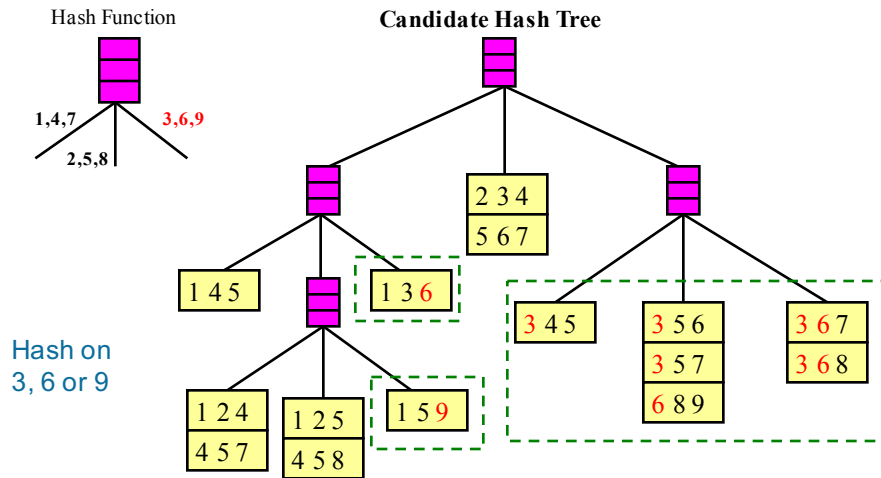


02/03/2018

Introduction to Data Mining

34

Support Counting Using a Hash Tree

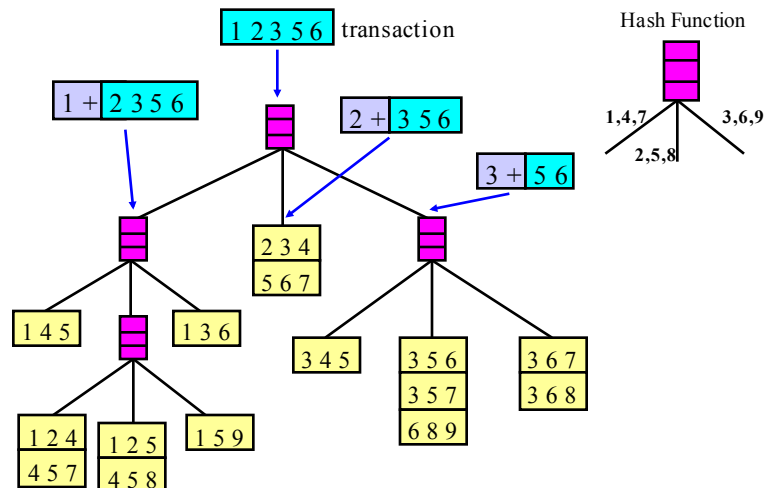


02/03/2018

Introduction to Data Mining

35

Support Counting Using a Hash Tree

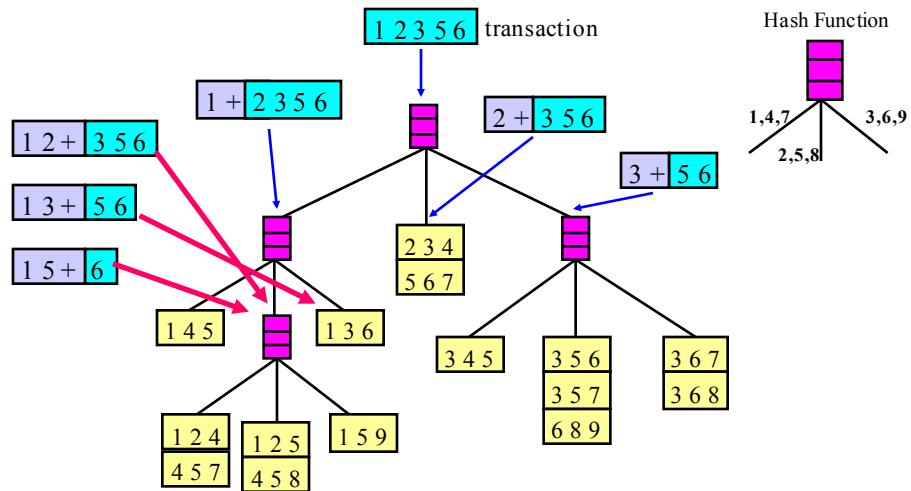


02/03/2018

Introduction to Data Mining

36

Support Counting Using a Hash Tree

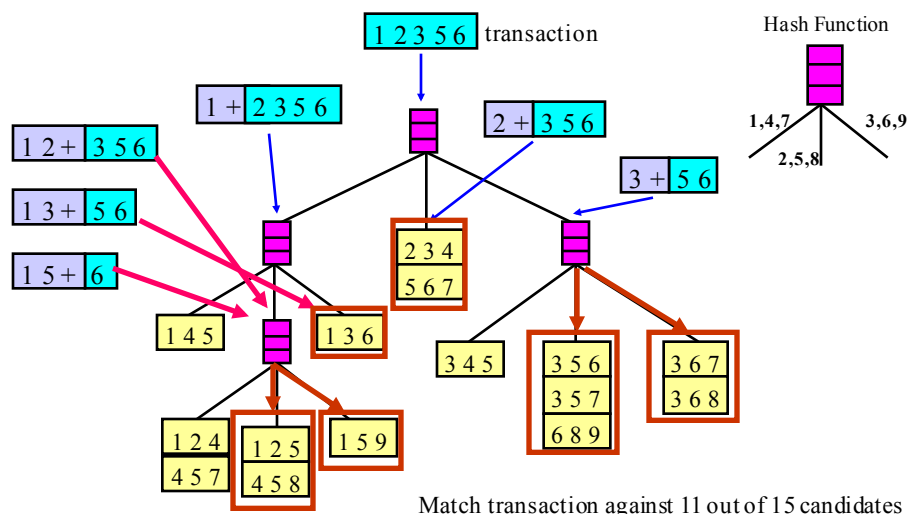


02/03/2018

Introduction to Data Mining

37

Support Counting Using a Hash Tree



Match transaction against 11 out of 15 candidates

02/03/2018

Introduction to Data Mining

38

Rule Generation

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

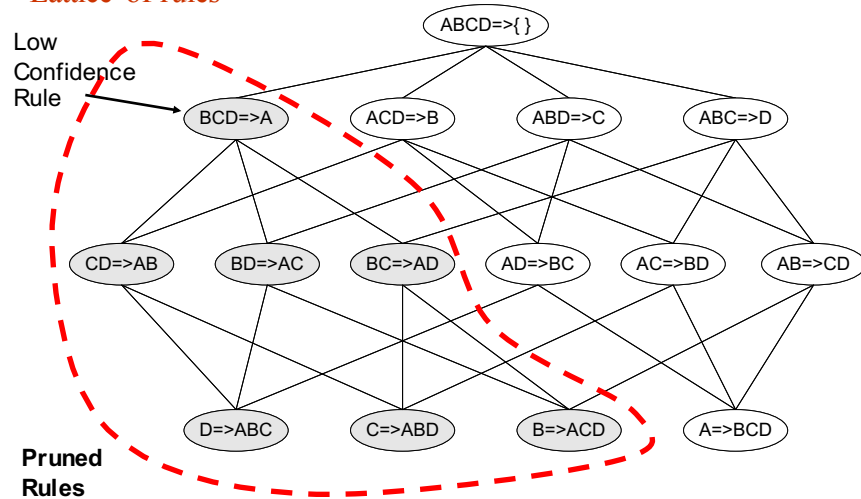
$ABC \rightarrow D,$	$ABD \rightarrow C,$	$ACD \rightarrow B,$	$BCD \rightarrow A,$
$A \rightarrow BCD,$	$B \rightarrow ACD,$	$C \rightarrow ABD,$	$D \rightarrow ABC$
$AB \rightarrow CD,$	$AC \rightarrow BD,$	$AD \rightarrow BC,$	$BC \rightarrow AD,$
$BD \rightarrow AC,$	$CD \rightarrow AB,$		
- If $|L| = k$, then there are $2^k - 2$ candidate association rules (ignoring $L \rightarrow \emptyset$ and $\emptyset \rightarrow L$)

Rule Generation

- In general, confidence does not have an anti-monotone property
 $c(ABC \rightarrow D)$ can be larger or smaller than $c(AB \rightarrow D)$
- 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:
$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$
 - Confidence is anti-monotone w.r.t. number of items on the RHS of the rule

Rule Generation for Apriori Algorithm

Lattice of rules



02/03/2018

Introduction to Data Mining

41

Association Analysis: Basic Concepts and Algorithms

Algorithms and Complexity

Factors Affecting Complexity of Apriori

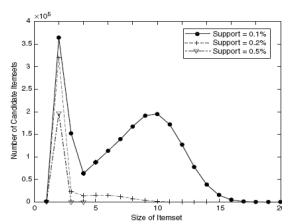
- Choice of minimum support threshold
 - lowering support threshold results in more frequent itemsets
 - this may increase number of candidates and max length of frequent itemsets
- Dimensionality (number of items) of the data set
 - more space is needed to store support count of each item
 - if number of frequent items also increases, both computation and I/O costs may also increase
- Size of database
 - since Apriori makes multiple passes, run time of algorithm may increase with number of transactions
- Average transaction width
 - transaction width increases with denser datasets
 - This may increase max length of frequent itemsets and traversals of hash tree (number of subsets in a transaction increases with its width)

02/03/2018

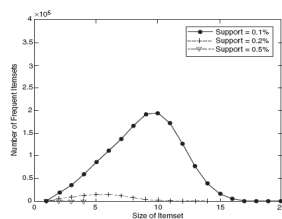
Introduction to Data Mining

43

Factors Affecting Complexity of Apriori

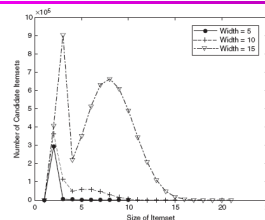


(a) Number of candidate itemsets.

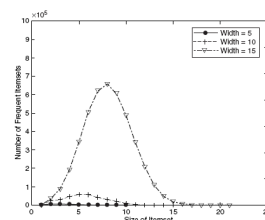


(b) Number of frequent itemsets.

Figure 6.13. Effect of support threshold on the number of candidate and frequent itemsets.



(a) Number of candidate itemsets.



(b) Number of frequent itemsets.

Figure 6.14. Effect of average transaction width on the number of candidate and frequent itemsets.

02/03/2018

Introduction to Data Mining

44

Compact Representation of Frequent Itemsets

- Some itemsets are redundant because they have identical support as their supersets

TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1

- Number of frequent itemsets = $3 \times \sum_{k=1}^{10} \binom{10}{k}$
- Need a compact representation

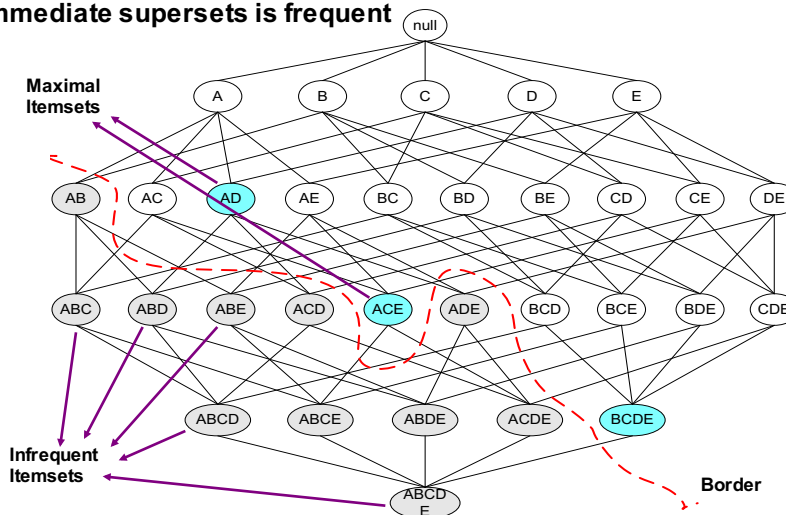
02/03/2018

Introduction to Data Mining

45

Maximal Frequent Itemset

An itemset is maximal frequent if it is frequent and none of its immediate supersets is frequent



02/03/2018

Introduction to Data Mining

46

What are the Maximal Frequent Itemsets in this Data?

TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1

Minimum support threshold = 5

02/03/2018

Introduction to Data Mining

47

An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

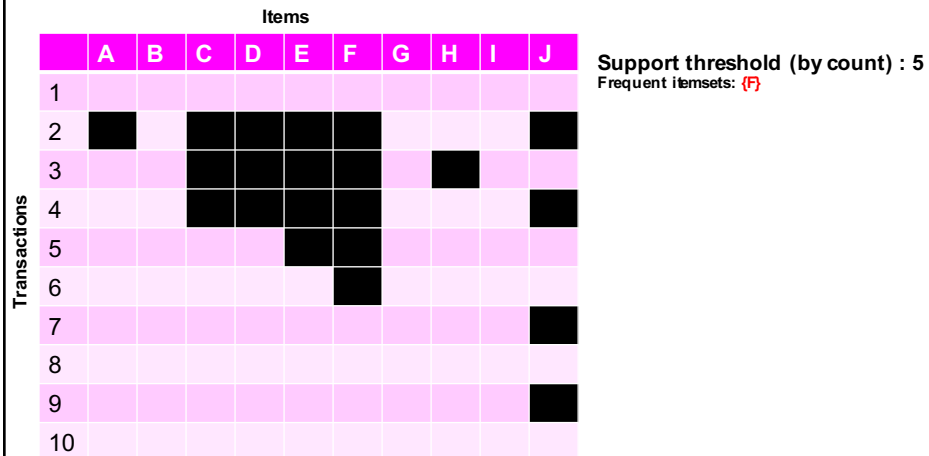
Support threshold (by count) : 5
Frequent itemsets: ?

02/03/2018

Introduction to Data Mining

48

An illustrative example

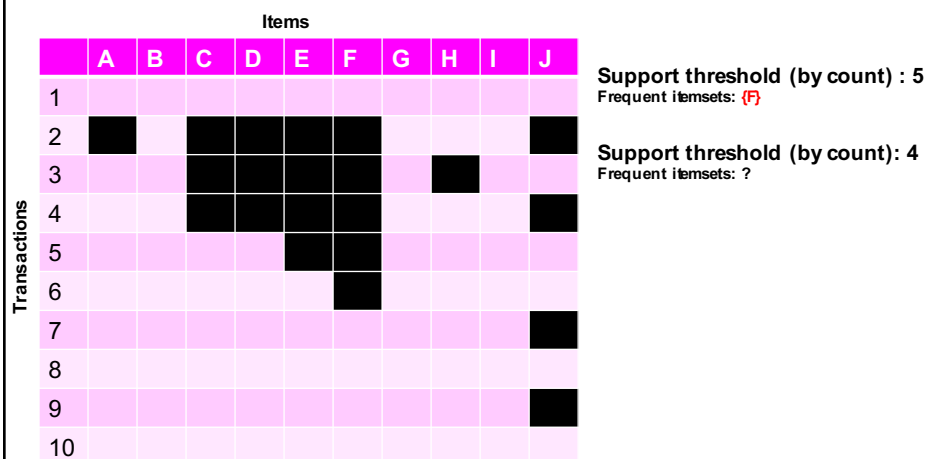


02/03/2018

Introduction to Data Mining

49

An illustrative example

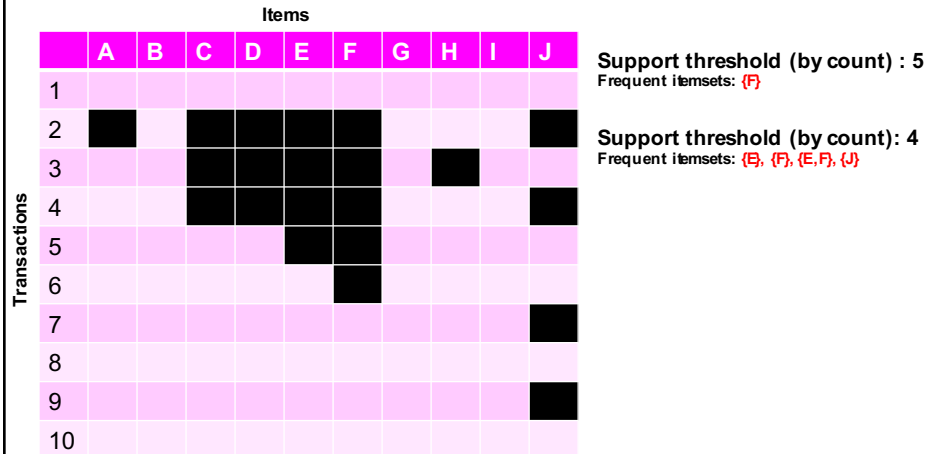


02/03/2018

Introduction to Data Mining

50

An illustrative example

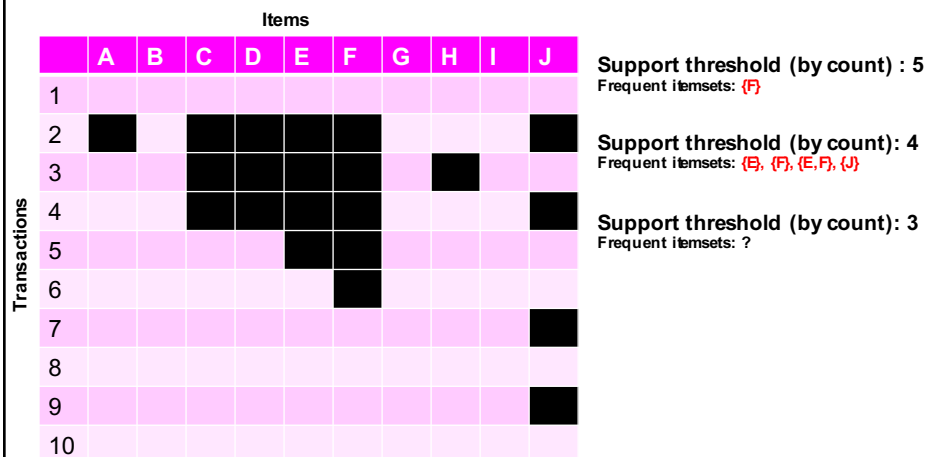


02/03/2018

Introduction to Data Mining

51

An illustrative example



02/03/2018

Introduction to Data Mining

52

An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Support threshold (by count) : 5
Frequent itemsets: {F}

Support threshold (by count): 4
Frequent itemsets: {E}, {F}, {E,F}, {J}

Support threshold (by count): 3
Frequent itemsets:
All subsets of {C,D,E,F} + {J}

An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Support threshold (by count) : 5
Frequent itemsets: {F}
Maximal itemsets: ?

Support threshold (by count): 4
Frequent itemsets: {E}, {F}, {E,F}, {J}
Maximal itemsets: ?

Support threshold (by count): 3
Frequent itemsets:
All subsets of {C,D,E,F} + {J}
Maximal itemsets: ?

An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Support threshold (by count) : 5
Frequent itemsets: {F}
Maximal itemsets: {F}

Support threshold (by count): 4
Frequent itemsets: {E}, {F}, {E,F}, {J}
Maximal itemsets: ?

Support threshold (by count): 3
Frequent itemsets:
All subsets of {C,D,E,F} + {J}
Maximal itemsets: ?

An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Support threshold (by count) : 5
Frequent itemsets: {F}
Maximal itemsets: {F}

Support threshold (by count): 4
Frequent itemsets: {E}, {F}, {E,F}, {J}
Maximal itemsets: {E,F}, {J}

Support threshold (by count): 3
Frequent itemsets:
All subsets of {C,D,E,F} + {J}
Maximal itemsets: ?

An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Support threshold (by count) : 5
Frequent itemsets: {F}
Maximal itemsets: {F}

Support threshold (by count): 4
Frequent itemsets: {E}, {F}, {E,F}, {J}
Maximal itemsets: {E,F}, {J}

Support threshold (by count): 3
Frequent itemsets:
All subsets of {C,D,E,F} + {J}
Maximal itemsets:
{C,D,E,F}, {J}

Another illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Support threshold (by count) : 5
Maximal itemsets: {A}, {B}, {C}

Support threshold (by count): 4
Maximal itemsets: {A,B}, {A,C}, {B,C}

Support threshold (by count): 3
Maximal itemsets: {A,B,C}

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

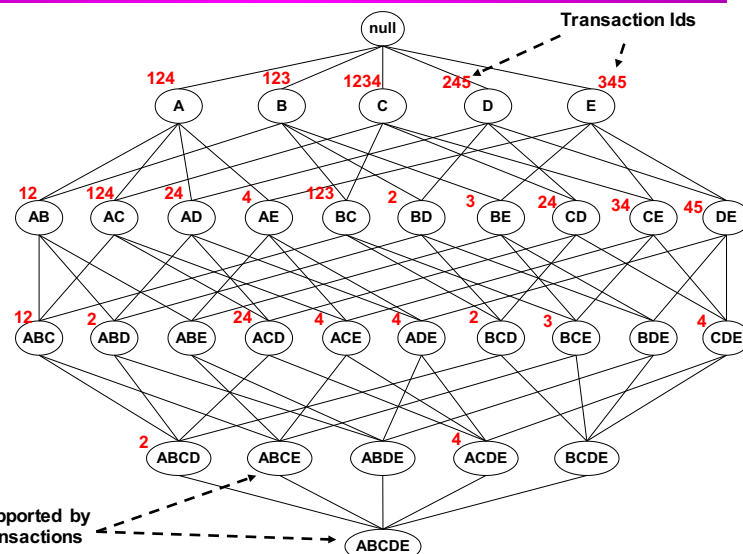
02/03/2018

Introduction to Data Mining

59

Maximal vs Closed Itemsets

TID	Items
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE



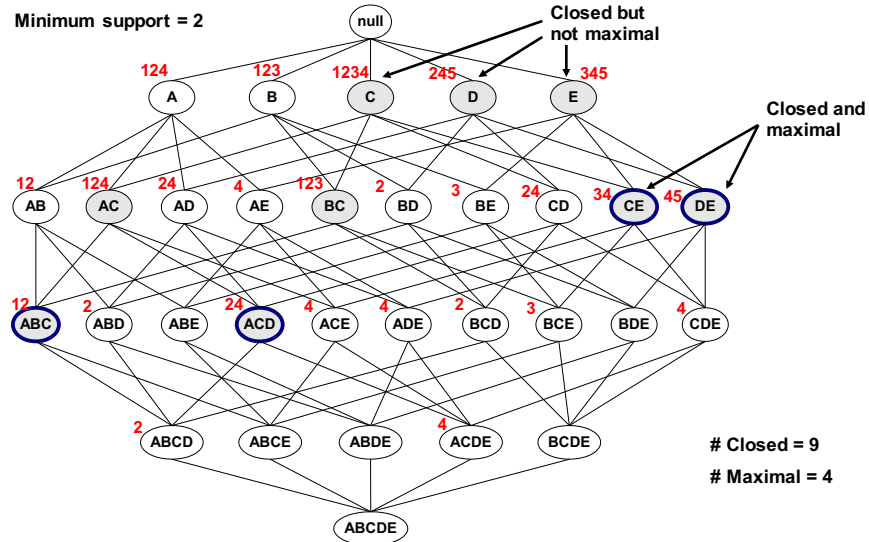
02/03/2018

Introduction to Data Mining

60

Maximal vs Closed Frequent Itemsets

Minimum support = 2



02/03/2018

Introduction to Data Mining

61

What are the Closed Itemsets in this Data?

TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1

02/03/2018

Introduction to Data Mining

62

Example 1

Transactions	Items										
	A	B	C	D	E	F	G	H	I	J	
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
9											
10											

Itemsets	Support (counts)	Closed itemsets
{C}	3	
{D}	2	
{C,D}	2	

Itemsets	Support (counts)	Closed itemsets
{C}	3	
{D}	2	
{C,D}	2	

02/03/2018

Introduction to Data Mining

63

Example 1

Transactions	Items										
	A	B	C	D	E	F	G	H	I	J	
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
9											
10											

Itemsets	Support (counts)	Closed itemsets
{C}	3	✓
{D}	2	
{C,D}	2	✓

Itemsets	Support (counts)	Closed itemsets
{C}	3	✓
{D}	2	
{C,D}	2	✓

02/03/2018

Introduction to Data Mining

64

Example 2

Transactions	Items										Itemsets	Support (counts)	Closed itemsets
	A	B	C	D	E	F	G	H	I	J			
1													
2											{C}	3	
3											{D}	2	
4											{E}	2	
5											{C,D}	2	
6											{C,E}	2	
7											{D,E}	2	
8											{C,D,E}	2	
9													
10													

02/03/2018

Introduction to Data Mining

65

Example 2

Transactions	Items										Itemsets	Support (counts)	Closed itemsets
	A	B	C	D	E	F	G	H	I	J			
1													
2											{C}	3	✓
3											{D}	2	
4											{E}	2	
5											{C,D}	2	
6											{C,E}	2	
7											{D,E}	2	
8											{C,D,E}	2	✓
9													
10													

02/03/2018

Introduction to Data Mining

66

Example 3

		Items										Transactions	
		A	B	C	D	E	F	G	H	I	J		
1													
2													
3													
4													
5													
6													
7													
8													
9													
10													

Closed itemsets: {C,D,E,F}, {C,F}

02/03/2018

Introduction to Data Mining

67

Example 4

		Items										Transactions	
		A	B	C	D	E	F	G	H	I	J		
1													
2													
3													
4													
5													
6													
7													
8													
9													
10													

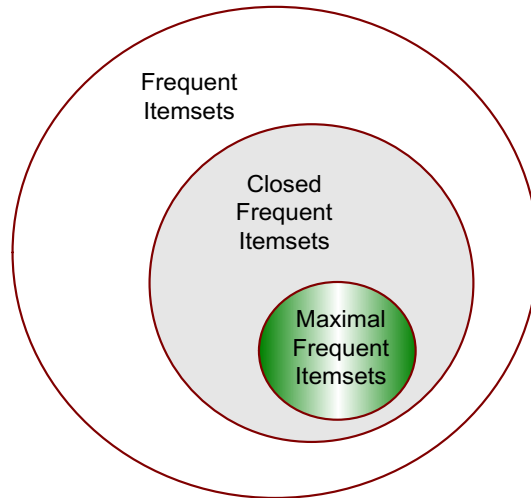
Closed itemsets: {C,D,E,F}, {C}, {F}

02/03/2018

Introduction to Data Mining

68

Maximal vs Closed Itemsets



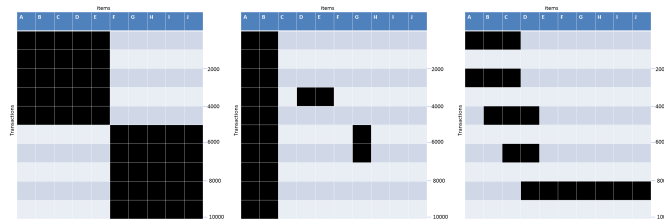
02/03/2018

Introduction to Data Mining

69

Example question

- Given the following transaction data sets (dark cells indicate presence of an item in a transaction) and a support threshold of 20%, answer the following questions



- What is the number of frequent itemsets for each dataset? Which dataset will produce the most number of frequent itemsets?
- Which dataset will produce the longest frequent itemset?
- Which dataset will produce frequent itemsets with highest maximum support?
- Which dataset will produce frequent itemsets containing items with widely varying support levels (i.e., itemsets containing items with mixed support, ranging from 20% to more than 70%)?
- What is the number of maximal frequent itemsets for each dataset? Which dataset will produce the most number of maximal frequent itemsets?
- What is the number of closed frequent itemsets for each dataset? Which dataset will produce the most number of closed frequent itemsets?

02/03/2018

Introduction to Data Mining

70

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	\bar{Y}	
X	f_{11}	f_{10}	f_{1+}
\bar{X}	f_{01}	f_{00}	f_{0+}
	f_{+1}	f_{+0}	N

f_{11} : support of X and Y
 f_{10} : support of X and \bar{Y}
 f_{01} : support of \bar{X} and Y
 f_{00} : support of \bar{X} and \bar{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	<u>Coffee</u>	
Tea	15	5	20
<u>Tea</u>	75	5	80
	90	10	100

Association Rule: Tea → Coffee

$$\text{Confidence} \equiv P(\text{Coffee}|\text{Tea}) = 15/20 = 0.75$$

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

So rule seems reasonable

Drawback of Confidence

	Coffee	<u>Coffee</u>	
Tea	15	5	20
<u>Tea</u>	75	5	80
	90	10	100

Association Rule: Tea → Coffee

$$\text{Confidence} = P(\text{Coffee}|\text{Tea}) = 15/20 = 0.75$$

but $P(\text{Coffee}) = 0.9$, which means knowing that a person drinks tea reduces the probability that the person drinks coffee!

$$\Rightarrow \text{Note that } P(\text{Coffee}|\text{Tea}) = 75/80 = 0.9375$$

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 Independence

- The criterion
confidence($X \rightarrow Y$) = support(Y)

is equivalent to:

- $P(Y|X) = P(Y)$
- $P(X,Y) = P(X) \times P(Y)$

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

$$\begin{aligned}
 \text{Lift} &= \frac{P(Y | X)}{P(Y)} \\
 \text{Interest} &= \frac{P(X, Y)}{P(X)P(Y)} \\
 PS &= P(X, Y) - P(X)P(Y) \\
 \phi - \text{coefficient} &= \frac{P(X, Y) - P(X)P(Y)}{\sqrt{P(X)[1 - P(X)]P(Y)[1 - P(Y)]}}
 \end{aligned}$$

lift is used for rules while
interest is used for itemsets

Example: Lift/Interest

	Coffee	<u>Coffee</u>	
Tea	15	5	20
<u>Tea</u>	75	5	80
	90	10	100

Association Rule: Tea → Coffee

Confidence = $P(\text{Coffee} | \text{Tea}) = 0.75$

but $P(\text{Coffee}) = 0.9$

⇒ Lift = $0.75/0.9 = 0.8333$ (< 1, therefore is negatively associated)

So, is it enough to use confidence/lift for pruning?

Lift or Interest

	Y	\bar{Y}	
X	10	0	10
\bar{X}	0	90	90
	10	90	100

$$Lift = \frac{0.1}{(0.1)(0.1)} = 10$$

	Y	\bar{Y}	
X	90	0	90
\bar{X}	0	10	10
	90	10	100

$$Lift = \frac{0.9}{(0.9)(0.9)} = 1.11$$

Statistical independence:

If $P(X,Y)=P(X)P(Y) \Rightarrow Lift = 1$

02/03/2018

Introduction to Data Mining

79

There are lots of measures proposed in the literature

#	Measure	Formula
1	ϕ -coefficient	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
2	Goodman-Kruskal's (λ)	$\frac{\sum_j \max_k P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
3	Odds ratio (α)	$\frac{P(A,B)P(\bar{A},\bar{B})}{P(A,\bar{B})P(\bar{A},B)}$
4	Yule's Q	$\frac{P(A,B)P(\bar{A},\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A},\bar{B}) + P(A,\bar{B})P(\bar{A},B)} = \frac{\alpha - 1}{\alpha + 1}$
5	Yule's Y	$\frac{\sqrt{P(A,B)P(\bar{A},\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A},\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}} = \frac{\sqrt{\alpha} - 1}{\sqrt{\alpha} + 1}$
6	Kappa (κ)	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
7	Mutual Information (M)	$\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}$
8	J-Measure (J)	$\min(-\sum_i P(A_i) \log P(A_i), -\sum_j P(B_j) \log P(B_j))$
9	Gini index (G)	$\max \left(P(A,B) \log \left(\frac{P(A B)}{P(B)} \right) + P(\bar{A},\bar{B}) \log \left(\frac{P(\bar{A} \bar{B})}{P(\bar{A})} \right), \right. \\ \left. P(A,B) \log \left(\frac{P(A B)}{P(B)} \right) + P(\bar{A},B) \log \left(\frac{P(\bar{A} B)}{P(\bar{A})} \right) \right)$
10	Support (s)	$P(A,B)$
11	Confidence (c)	$\max\{P(B A), P(A B)\}$
12	Laplace (L)	$\max \left(\frac{NP(A,B)+1}{NP(A)+1}, \frac{NP(A,B)+1}{NP(B)+1} \right)$
13	Conviction (V)	$\max \left(\frac{P(A)P(\bar{B})}{P(A,B)}, \frac{P(B)P(\bar{A})}{P(A,B)} \right)$
14	Interest (I)	$\frac{P(A,B)}{P(A)P(B)}$
15	cosine (IS)	$\frac{P(A,B)}{\sqrt{P(A)P(B)}}$
16	Piatetsky-Shapiro's (PS)	$P(A,B) - P(A)P(B)$
17	Certainty factor (F)	$\max \left(\frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)} \right)$
18	Added Value (AV)	$\max\{P(B A) - P(B), P(A B) - P(A)\}$
19	Collective strength (S)	$\frac{P(A,B) + P(\bar{A},\bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A,B) - P(\bar{A},\bar{B})}$
20	Jaccard (ζ)	$\frac{P(A,B)}{P(A) + P(B) - P(A,B)}$
21	Kloggen (K)	$\sqrt{P(A,B)} \max\{P(B A) - P(B), P(A B) - P(A)\}$

02/03/2018

Comparing Different Measures

10 examples of contingency tables:

Example	f_{11}	f_{10}	f_{01}	f_{00}
E1	8123	83	424	1370
E2	8330	2	622	1046
E3	9481	94	127	298
E4	3954	3080	5	2961
E5	2886	1363	1320	4431
E6	1500	2000	500	6000
E7	4000	2000	1000	3000
E8	4000	2000	2000	2000
E9	1720	7121	5	1154
E10	61	2483	4	7452

Rankings of contingency tables using various measures:

#	ϕ	λ	α	Q	Y	κ	M	J	G	s	c	L	V	I	IS	PS	F	AV	S	ζ	K
E1	1	1	3	3	3	1	2	2	1	3	5	5	4	6	2	2	4	6	1	2	5
E2	2	2	1	1	1	2	1	3	2	2	1	1	1	8	3	5	1	8	2	3	6
E3	3	3	4	4	4	3	3	8	7	1	4	4	6	10	1	8	6	10	3	1	10
E4	4	7	2	2	2	5	4	1	3	6	2	2	2	4	4	1	2	3	4	5	1
E5	5	4	8	8	8	4	7	5	4	7	9	9	9	3	6	3	9	4	5	6	3
E6	6	6	7	7	7	7	6	4	6	9	8	8	7	2	8	6	7	2	7	8	2
E7	7	5	9	9	9	6	8	6	5	4	7	7	8	5	5	4	8	5	6	4	4
E8	8	9	10	10	10	8	10	10	8	4	10	10	10	9	7	7	10	9	8	7	9
E9	9	9	5	5	5	9	9	7	9	8	3	3	3	7	9	9	3	7	9	9	8
E10	10	8	6	6	6	10	5	9	10	10	6	6	5	1	10	10	5	1	10	10	7

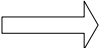
02/03/2018

Introduction to Data Mining

81

Property under Variable Permutation

	B	$\overline{\mathbf{B}}$
A	p	q
$\overline{\mathbf{A}}$	r	s



	A	$\overline{\mathbf{A}}$
B	p	r
$\overline{\mathbf{B}}$	q	s

Does $M(A,B) = M(B,A)$?

Symmetric measures:

- ◆ support, lift, collective strength, cosine, Jaccard, etc

Asymmetric measures:

- ◆ confidence, conviction, Laplace, J-measure, etc

02/03/2018

Introduction to Data Mining

82

Property under Row/Column Scaling

Grade-Gender Example (Mosteller, 1968):

	Female	Male	
High	2	3	5
Low	1	4	5
	3	7	10

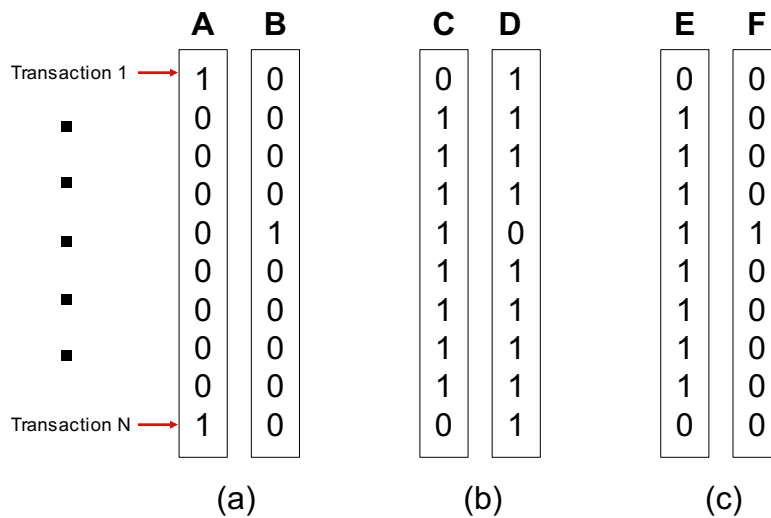
	Female	Male	
High	4	30	34
Low	2	40	42
	6	70	76

↓
2x ↓
10x

Mosteller:

Underlying association should be independent of the relative number of male and female students in the samples

Property under Inversion Operation



Example: ϕ -Coefficient

- ϕ -coefficient is analogous to correlation coefficient for continuous variables

	Y	\bar{Y}	
X	60	10	70
\bar{X}	10	20	30
	70	30	100

	Y	\bar{Y}	
X	20	10	30
\bar{X}	10	60	70
	30	70	100

$$\phi = \frac{0.6 - 0.7 \times 0.7}{\sqrt{0.7 \times 0.3 \times 0.7 \times 0.3}} = 0.5238$$

$$\phi = \frac{0.2 - 0.3 \times 0.3}{\sqrt{0.7 \times 0.3 \times 0.7 \times 0.3}} = 0.5238$$

ϕ Coefficient is the same for both tables

Property under Null Addition

	B	\bar{B}
A	p	q
\bar{A}	r	s

→

	B	\bar{B}
A	p	q
\bar{A}	r	s + k

Invariant measures:

- ◆ support, cosine, Jaccard, etc

Non-invariant measures:

- ◆ correlation, Gini, mutual information, odds ratio, etc

Different Measures have Different Properties

Symbol	Measure	Inversion	Null Addition	Scaling
ϕ	ϕ -coefficient	Yes	No	No
α	odds ratio	Yes	No	Yes
κ	Cohen's	Yes	No	No
I	Interest	No	No	No
IS	Cosine	No	Yes	No
PS	Piatetsky-Shapiro's	Yes	No	No
S	Collective strength	Yes	No	No
ζ	Jaccard	No	Yes	No
h	All-confidence	No	No	No
s	Support	No	No	No

Simpson's Paradox

Buy HDTV	Buy Exercise Machine		
	Yes	No	
Yes	99	81	180
No	54	66	120
	153	147	300

$$c(\{\text{HDTV} = \text{Yes}\} \rightarrow \{\text{Exercise Machine} = \text{Yes}\}) = 99/180 = 55\%$$

$$c(\{\text{HDTV} = \text{No}\} \rightarrow \{\text{Exercise Machine} = \text{Yes}\}) = 54/120 = 45\%$$

=> Customers who buy HDTV are more likely to buy exercise machines

Simpson's Paradox

Customer Group	Buy HDTV	Buy Exercise Machine		Total
		Yes	No	
College Students	Yes	1	9	10
	No	4	30	34
Working Adult	Yes	98	72	170
	No	50	36	86

College students:

$$c(\{\text{HDTV} = \text{Yes}\} \rightarrow \{\text{Exercise Machine} = \text{Yes}\}) = 1/10 = 10\%$$

$$c(\{\text{HDTV} = \text{No}\} \rightarrow \{\text{Exercise Machine} = \text{Yes}\}) = 4/34 = 11.8\%$$

Working adults:

$$c(\{\text{HDTV} = \text{Yes}\} \rightarrow \{\text{Exercise Machine} = \text{Yes}\}) = 98/170 = 57.7\%$$

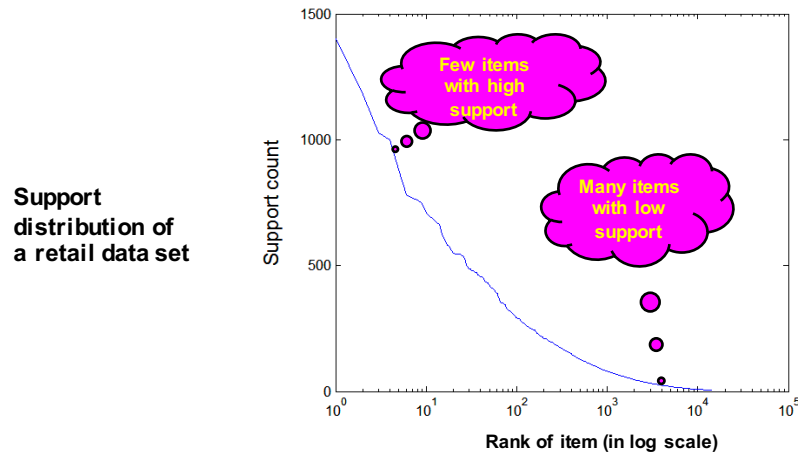
$$c(\{\text{HDTV} = \text{No}\} \rightarrow \{\text{Exercise Machine} = \text{Yes}\}) = 50/86 = 58.1\%$$

Simpson's Paradox

- Observed relationship in data may be influenced by the presence of other confounding factors (hidden variables)
 - Hidden variables may cause the observed relationship to disappear or reverse its direction!
- Proper stratification is needed to avoid generating spurious patterns

Effect of Support Distribution on Association Mining

- Many real data sets have skewed support distribution



02/03/2018

Introduction to Data Mining

91

Effect of Support Distribution

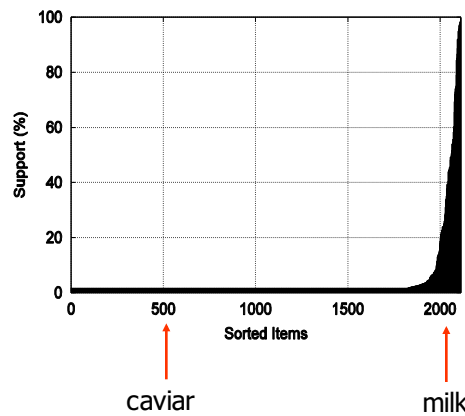
- Difficult to set the appropriate *minsup* threshold
 - If *minsup* is too high, we could miss itemsets involving interesting rare items (e.g., {caviar, vodka})
 - If *minsup* is too low, it is computationally expensive and the number of itemsets is very large

02/03/2018

Introduction to Data Mining

92

Cross-Support Patterns



A cross-support pattern involves items with varying degree of support

- Example: {caviar,milk}

How to avoid such patterns?

A Measure of Cross Support

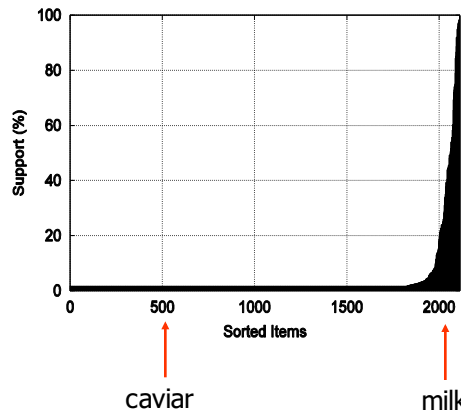
- Given an itemset, $X = \{x_1, x_2, \dots, x_d\}$, with d items, we can define a measure of cross support, r , for the itemset

$$r(X) = \frac{\min\{s(x_1), s(x_2), \dots, s(x_d)\}}{\max\{s(x_1), s(x_2), \dots, s(x_d)\}}$$

where $s(x_i)$ is the support of item x_i

- Can use $r(X)$ to prune cross support patterns, but not to avoid them

Confidence and Cross-Support Patterns



Observation:

$\text{conf}(\text{caviar} \rightarrow \text{milk})$ is very high
but
 $\text{conf}(\text{milk} \rightarrow \text{caviar})$ is very low

Therefore,

$\min(\text{conf}(\text{caviar} \rightarrow \text{milk}), \text{conf}(\text{milk} \rightarrow \text{caviar}))$

is also very low

H-Confidence

- To avoid patterns whose items have very different support, define a new evaluation measure for itemsets
 - Known as **h-confidence** or **all-confidence**
- Specifically, given an itemset $X = \{x_1, x_2, \dots, x_d\}$
 - h-confidence is the minimum confidence of any association rule formed from itemset X
 - $\text{hconf}(X) = \min(\text{conf}(X_1 \rightarrow X_2))$,
where $X_1, X_2 \subset X, X_1 \cap X_2 = \emptyset, X_1 \cup X_2 = X$
For example: $X_1 = \{x_1, x_2\}, X_2 = \{x_3, \dots, x_d\}$

H-Confidence ...

- But, given an itemset $X = \{x_1, x_2, \dots, x_d\}$
 - What is the lowest confidence rule you can obtain from X ?
 - Recall $\text{conf}(X_1 \rightarrow X_2) = s(X_1 \cup X_2) / \text{support}(X_1)$
 - ◆ The numerator is fixed: $s(X_1 \cup X_2) = s(X)$
 - ◆ Thus, to find the lowest confidence rule, we need to find the X_1 with highest support
 - ◆ Consider only rules where X_1 is a single item, i.e., $\{x_1\} \rightarrow X - \{x_1\}$, $\{x_2\} \rightarrow X - \{x_2\}$, ..., or $\{x_d\} \rightarrow X - \{x_d\}$

$$\begin{aligned} \text{hconf}(X) &= \min \left\{ \frac{s(X)}{s(x_1)}, \frac{s(X)}{s(x_2)}, \dots, \frac{s(X)}{s(x_d)} \right\} \\ &= \frac{s(X)}{\max\{s(x_1), s(x_2), \dots, s(x_d)\}} \end{aligned}$$

Cross Support and H-confidence

- By the anti-montone property of support

$$s(X) \leq \min\{s(x_1), s(x_2), \dots, s(x_d)\}$$

- Therefore, we can derive a relationship between the h-confidence and cross support of an itemset

$$\begin{aligned} \text{hconf}(X) &= \frac{s(X)}{\max\{s(x_1), s(x_2), \dots, s(x_d)\}} \\ &\leq \frac{\min\{s(x_1), s(x_2), \dots, s(x_d)\}}{\max\{s(x_1), s(x_2), \dots, s(x_d)\}} \\ &= r(X) \end{aligned}$$

Thus, $\text{hconf}(X) \leq r(X)$

Cross Support and H-confidence ...

- Since, $\text{hconf}(X) \leq r(X)$, we can eliminate cross support patterns by finding patterns with $\text{h-confidence} < h_c$, a user set threshold
- Notice that

$$0 \leq \text{hconf}(X) \leq r(X) \leq 1$$

- Any itemset satisfying a given h-confidence threshold, h_c , is called a **hyperclique**
- H-confidence can be used instead of or in conjunction with support

Properties of Hypercliques

- Hypercliques are itemsets, but not necessarily frequent itemsets
 - Good for finding low support patterns
- H-confidence is anti-monotone
- Can define closed and maximal hypercliques in terms of h-confidence
 - A hyperclique X is closed if none of its immediate supersets has the same h-confidence as X
 - A hyperclique X is maximal if $\text{hconf}(X) \leq h_c$ and none of its immediate supersets, Y , have $\text{hconf}(Y) \leq h_c$

Properties of Hypercliques ...

- Hypercliques have the high-affinity property
 - Think of the individual items as sparse binary vectors
 - h-confidence gives us information about their pairwise Jaccard and cosine similarity
 - ◆ Assume x_1 and x_2 are any two items in an itemset X
 - ◆ $\text{Jaccard}(x_1, x_2) \geq \text{hconf}(X)/2$
 - ◆ $\cos(x_1, x_2) \geq \text{hconf}(X)$
 - Hypercliques that have a high h-confidence consist of very similar items as measured by Jaccard and cosine
- The items in a hyperclique cannot have widely different support
 - Allows for more efficient pruning

02/03/2018

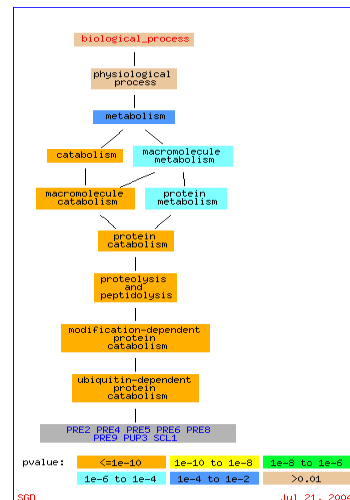
Introduction to Data Mining

101

Example Applications of Hypercliques

- Hypercliques are used to find strongly coherent groups of items
 - Words that occur together in documents
 - Proteins in a protein interaction network

In the figure at the right, a gene ontology hierarchy for biological process shows that the identified proteins in the hyperclique (PRE2, ..., SCL1) perform the same function and are involved in the same biological process



02/03/2018

Introduction to Data Mining

102