# Introduction to Data Mining

## Chapter 4
## Association Analysis:
## Basic Concepts and Algorithms

*Source: revised from slides provided by Tan, Steinbach, Karpatne, and 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**

{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 ($\sigma$)**
  - 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 $X \rightarrow Y$, where X and Y are itemsets

  - Example:
    {Milk, Diaper} → {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:

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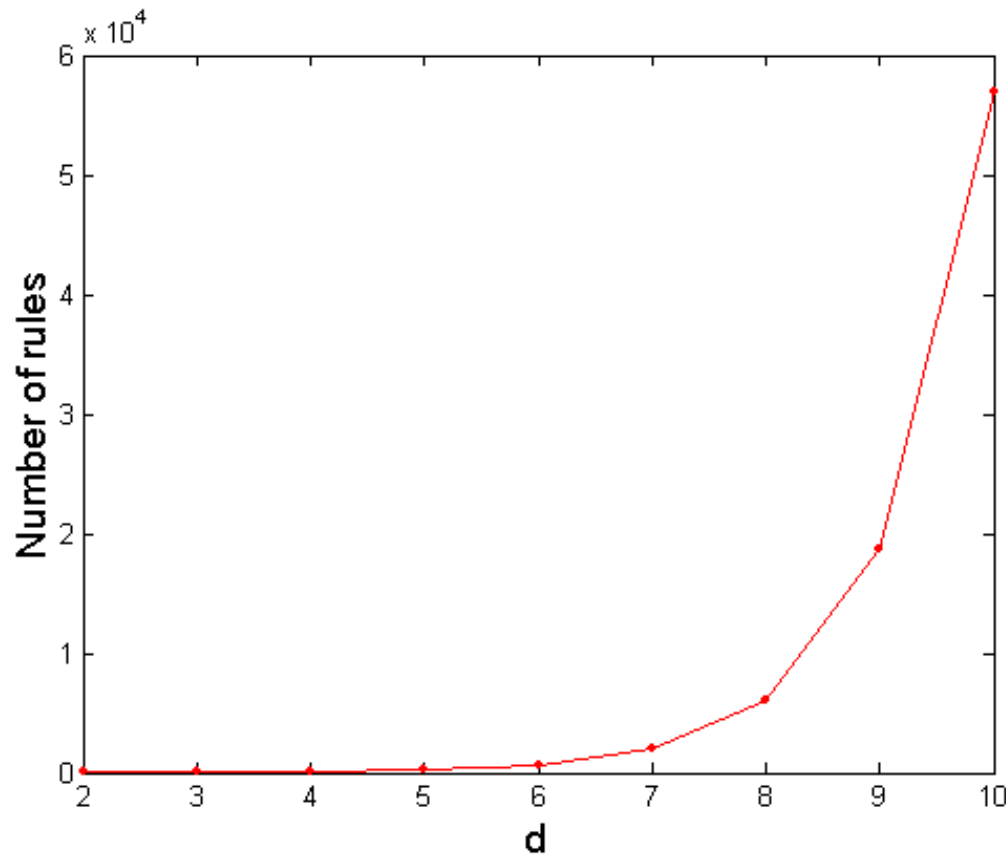
$$s = \frac{\sigma(Milk, Diaper, Beer)}{|T|} = \frac{2}{5} = 0.4$$

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

# 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!

# Computational Complexity

- Given d unique items:
  - Total number of itemsets = $2^d - 1$
  - 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**

# 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:

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

## 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 $\geq$ 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

# Frequent Itemset Generation

**Itemset lattice**



**Given d items, there are ($2^d − 1$) possible candidate itemsets**

# 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**

| 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 |

N

w

**List of Candidates**

M

– Match each transaction against every candidate

– Complexity ~ O(NMw) => Expensive since M = $2^d$ !!!

# 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
  - If the width of a transaction is small (e.g., 2 items), this transaction can be removed before searching for longer itemsets  (e.g., itemsets of size 3 and larger)

- 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



Found to be Infrequent

Pruned supersets

# 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 |

**6 distinct items**

**Candidate: 6**

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,

$$\binom{6}{1} + \binom{6}{2} + \binom{6}{3}$$
$$6 + 15 + 20 = 41$$

# 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 |

**Frequent: 4**

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,
$$\binom{6}{1} + \binom{6}{2} + \binom{6}{3}$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
$$\binom{6}{1} + \binom{4}{2} + \binom{4}{3}$$
$$6 + 6 + 4 = 16$$

# Illustrating Apriori Principle

Items (1-itemsets)

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

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

**Candidate: 6**

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

**Minimum Support = 3**

If every subset is considered,

$$\binom{6}{1} + \binom{6}{2} + \binom{6}{3}$$

$$6 + 15 + 20 = 41$$

With support-based pruning,

$$\binom{6}{1} + \binom{4}{2} + \binom{4}{3}$$

$$6 + 6 + 4 = 16$$

# Illustrating Apriori Principle

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

Items (1-itemsets)

Frequent: 4

Pairs (2-itemsets)

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

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
$$\binom{6}{1} + \binom{6}{2} + \binom{6}{3}$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
$$\binom{6}{1} + \binom{4}{2} + \binom{4}{3}$$
$$6 + 6 + 4 = 16$$

# Illustrating Apriori Principle

Items (1-itemsets)

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

Pairs (2-itemsets)

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

(No need to generate candidates involving Coke or Eggs)

Candidate: 4

Triplets (3-itemsets)

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

Candidate: 1

Minimum Support = 3

If every subset is considered,
$$\binom{6}{1} + \binom{6}{2} + \binom{6}{3}$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
$$6 + 6 + 4 = 16$$
$$6 + 6 + 1 = 13$$

# Illustrating Apriori Principle

Candidate: 6, frequent: 4

Items (1-itemsets)

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

Candidate: 6, frequent: 4

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

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

Candidate: 1, frequent: 0

Triplets (3-itemsets)

| Itemset | Count |
|---------|-------|
| {Bread, Diaper, Milk} | 2 |

Minimum Support = 3

If every subset is considered,
$$\binom{6}{1} + \binom{6}{2} + \binom{6}{3}$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
$$6 + 6 + 4 = 16$$
$$6 + 6 + 1 = 13$$

# Apriori Algorithm

- $F_k$: frequent k-itemsets
- $C_k$: candidate k-itemsets

- **Algorithm**
  - Let k=1
  - Generate $F_1$ = {frequent 1-itemsets}
  - Repeat until $F_k$ is empty
    - **Candidate Generation**: Generate $C_{k+1}$ from $F_k$
    - **Candidate Pruning**: Prune candidate itemsets in $C_{k+1}$ containing subsets of length k that are infrequent
    - **Support Counting**: Count the support of each candidate in $C_{k+1}$ by scanning the DB
    - **Candidate Elimination**: Eliminate candidates in $C_{k+1}$ that are infrequent, leaving only those that are frequent => $F_{k+1}$

# Candidate Generation: Brute-force method



**Figure 6.6.** A brute-force method for generating candidate 3-itemsets.

**Figure 6.7.** Generating and pruning candidate $k$-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: $F_{k-1}$ x $F_{k-1}$ 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(**AB**C, **AB**D) = **AB**CD
  - Merge(**AB**C, **AB**E) = **AB**CE
  - Merge(**AB**D, **AB**E) = **AB**DE

  - Do not merge(**A**BD,**A**CD) 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

- $C_4$ = {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: $C_4$ = {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(A**BC**, **BC**D) = A**BC**D
  - Merge(A**BD**, **BD**E) = A**BD**E
  - Merge(A**CD**, **CD**E) = A**CD**E
  - Merge(B**CD**, **CD**E) = B**CD**E

# 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

- $C_4$ = {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: $C_4$ = {ABCD}

# Illustrating Apriori Principle

Items (1-itemsets)

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

Pairs (2-itemsets)

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

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
$$\binom{6}{1} + \binom{6}{2} + \binom{6}{3}$$
$$6 + 15 + 20 = 41$$
With support-based pruning,
$$6 + 6 + 1 = 13$$

Triplets (3-itemsets)

| Itemset | Count |
|---------|-------|
| {Bread, Diaper, Milk} | 2 |

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.

# 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

**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 |

N

**Hash Structure**

k

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)?**

Transaction, t

| 1 2 3 5 6 |

*Level 1*

**1** | 2 3 5 6      **2** | 3 5 6      **3** | 5 6

*Level 2*

**1 2** | 3 5 6    **1 3** | 5 6    **1 5** | 6    **2 3** | 5 6    **2 5** | 6    **3 5** | 6

```
1 2 3        1 3 5                    2 3 5
1 2 5        1 3 6        1 5 6       2 3 6       2 5 6       3 5 6
1 2 6
```

*Level 3*        Subsets of 3 items

# 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)**

# Support Counting Using a Hash Tree

Hash Function

**Candidate Hash Tree**

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

Candidate 3-itemsets:

{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}

Hash on 1, 4 or 7

1 4 5    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

1 2 4
4 5 7

1 2 5
4 5 8

1 5 9

# Support Counting Using a Hash Tree

Hash Function

Candidate Hash Tree

**Candidate 3-itemsets:**

**{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}**

1,4,7          3,6,9

2,5,8

Hash on
2, 5 or 8

2 3 4
5 6 7

1 4 5          1 3 6

3 4 5          3 5 6          3 6 7
               3 5 7          3 6 8
               6 8 9

1 2 4          1 2 5          1 5 9
4 5 7          4 5 8

# Support Counting Using a Hash Tree



Hash Function

Candidate Hash Tree

**Candidate 3-itemsets:**

{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}

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

Hash on 3, 6 or 9

# Support Counting Using a Hash Tree

1 2 3 5 6 transaction

Hash Function

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

1 + 2 3 5 6

2 + 3 5 6

3 + 5 6

2 3 4
5 6 7

1 4 5

1 3 6

3 4 5

3 5 6
3 5 7
6 8 9

3 6 7
3 6 8

1 2 4
4 5 7

1 2 5
4 5 8

1 5 9

# Support Counting Using a Hash Tree

1 2 3 5 6  transaction

Hash Function

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

1 + 2 3 5 6

2 + 3 5 6

1 2 + 3 5 6

1 3 + 5 6

1 5 + 6

3 + 5 6

2 3 4
5 6 7

1 4 5

1 3 6

3 4 5

3 5 6
3 5 7
6 8 9

3 6 7
3 6 8

1 2 4
4 5 7

1 2 5
4 5 8

1 5 9

# Support Counting Using a Hash Tree

1 2 3 5 6  transaction

Hash Function

1 + 2 3 5 6

2 + 3 5 6

1 2 + 3 5 6

1,4,7      3,6,9

2,5,8

1 3 + 5 6

3 + 5 6

1 5 + 6

3 5 + 6

2 3 4
5 6 7

1 4 5

1 3 6

3 4 5

3 5 6
3 5 7
6 8 9

3 6 7
3 6 8

1 2 4
4 5 7

1 2 5
4 5 8

1 5 9

- 5 out of the 9 leaf nodes are visited
- 9 out of 15 candidates are compared against the transaction

# Some Efficient and Scalable Frequent Itemset Mining Methods

Data Mining: Concepts and Techniques, 2nd ed.

By Jiawei Han and Micheline Kamber

*Source: Revised from Jiawei Han and Micheline Kamber's slides*

# 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

  - Hash entries: {ab, ad, ae} {bd, be, de} …

  - Frequent 1-itemset: a, b, d, e

  - ab is not a candidate 2-itemset if the sum of count of {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*

$H_2$

Create hash table $H_2$
using hash function
$h(x, y) = ((order\ of\ x) \times 10 + (order\ of\ y))\ mod\ 7$

| bucket address | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| bucket count | 2 | 2 | 4 | 2 | 2 | 4 | 4 |
| bucket contents | {I1, I4}<br>{I3, I5} | {I1, I5}<br>{I1, I5} | {I2, I3}<br>{I2, I3}<br>{I2, I3}<br>{I2, I3} | {I2, I4}<br>{I2, I4} | {I2, I5}<br>{I2, I5} | {I1, I2}<br>{I1, I2}<br>{I1, I2}<br>{I1, I2} | {I1, I3}<br>{I1, I3}<br>{I1, I3}<br>{I1, I3} |

# CHARM: Mining by Exploring Vertical Data Format

- Vertical format: $t(AB) = \{T_{11}, T_{25}, ...\}$
  - tid-list: list of trans.-ids containing an itemset
- Deriving closed patterns based on vertical intersections
  - $t(X) = t(Y)$: X and Y always happen together
  - $t(X) \subset t(Y)$: transaction having X always has Y
- Using diffset to accelerate mining
  - Only keep track of differences of tids
  - $t(X) = \{T_1, T_2, T_3\}$, $t(XY) = \{T_1, T_3\}$
  - Diffset $(XY, X) = \{T_2\}$
- Eclat/MaxEclat (Zaki et al. @KDD'97), VIPER(P. Shenoy et al.@SIGMOD'00), CHARM (Zaki & Hsiao@SDM'02)

# Example

| itemset | TID_set |
|---------|---------|
| I1 | {T100, T400, T500, T700, T800, T900} |
| I2 | {T100, T200, T300, T400, T600, T800, T900} |
| I3 | {T300, T500, T600, T700, T800, T900} |
| I4 | {T200, T400} |
| I5 | {T100, T800} |

**minsup = 2**

| itemset | TID_set |
|---------|---------|
| {I1,I2} | {T100, T400, T800, T900} |
| {I1,I3} | {T500, T700, T800, T900} |
| {I1,I4} | {T400} |
| {I1,I5} | {T100, T800} |
| {I2,I3} | {T300, T600, T800, T900} |
| {I2,I4} | {T200, T400} |
| {I2,I5} | {T100, T800} |
| {I3,I5} | {T800} |

**Intersect each pair of $F_1$**

- **Generate $C_3$ from $F_2 \times F_2$ (similar to Apriori)**
- **Intersect corresponding TID_sets**

| itemset | TID_set |
|---------|---------|
| {I1,I2,I3} | {T800, T900} |
| {I1,I2,I5} | {T100, T800} |

# Bottleneck of Frequent-pattern Mining

- Multiple database scans are costly

- Mining long patterns needs many passes of scanning and generates lots of candidates

  - To find frequent itemset $i_1 i_2 \ldots i_{100}$

    - \# of scans: 100

    - \# of Candidates: $\binom{100}{1} + \binom{100}{2} + \ldots + \binom{100}{100} = 2^{100} - 1 = 1.27 * 10^{30}$ !

- Bottleneck: candidate-generation-and-test

- Can we avoid candidate generation?

# Mining Frequent Patterns Without Candidate Generation

- Grow long patterns from short ones using local frequent items

    - "abc" is a frequent pattern

    - Get all transactions having "abc": DB|abc

    - "d" is a local frequent item in DB|abc → abcd is a frequent pattern

# Construct FP-tree from a Transaction Database

| TID | Items bought | (ordered) frequent items |
|-----|--------------|--------------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o, w} | {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} |

*min_support = 3*

1. Scan DB once, find frequent 1-itemset (single item pattern)

2. Sort frequent items in frequency descending order, f-list

3. Scan DB again, construct FP-tree

**Header Table**

| Item | frequency | head |
|------|-----------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |
| p | 3 | |



F-list=f-c-a-b-m-p

# Benefits of the FP-tree Structure

- Completeness
  - Preserve complete information for frequent pattern mining
  - Never break a long pattern of any transaction
- Compactness
  - Reduce irrelevant info—infrequent items are gone
  - Items in frequency descending order: the more frequently occurring, the more likely to be shared
  - Never be larger than the original database (not count node-links and the *count* field)

# Partition Patterns and Databases

- Frequent patterns can be partitioned into subsets according to f-list
  - F-list=f-c-a-b-m-p
  - Patterns containing p
  - Patterns having m but no p
  - …
  - Patterns having c but no a nor b, m, p
  - Pattern f
- Completeness and non-redundancy

# 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

**Header Table**

| Item | frequency | head |
|------|-----------|------|
| f    | 4         |      |
| c    | 4         |      |
| a    | 3         |      |
| b    | 3         |      |
| m    | 3         |      |
| p    | 3         |      |

{}

f:4    c:1

c:3   b:1   b:1

a:3         p:1

m:2   b:1

p:2   m:1

*Conditional* **pattern bases**

| item | 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       |

# From Conditional Pattern-bases to Conditional FP-trees

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

**Header Table**

| Item | frequency | head |
|------|-----------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |
| p | 3 | |

```
        {}
      /    \
   f:4      c:1
    |        |
   c:3  b:1  b:1
    |        |
   a:3      p:1
   /  \
 m:2   b:1
  |     |
 p:2   m:1
```

*m-conditional* **pattern base:**
*fca:2, fcab:1*

➔

```
 {}
  |
 f:3
  |
 c:3
  |
 a:3
```
*m-conditional* **FP-tree**

➔

**All frequent patterns relate to *m***

*m,*

*fm, cm, am,*

*fcm, fam, cam,*

*fcam*

# Example



Support count

An FP-tree registers compressed, frequent pattern information.

Mining the FP-tree by creating conditional (sub-)pattern bases.

| Item | Conditional Pattern Base | Conditional FP-tree | Frequent Patterns Generated |
|------|--------------------------|---------------------|------------------------------|
| I5 | {{I2, I1: 1}, {I2, I1, I3: 1}} | ⟨I2: 2, I1: 2⟩ | {I2, I5: 2}, {I1, I5: 2}, {I2, I1, I5: 2} |
| I4 | {{I2, I1: 1}, {I2: 1}} | ⟨I2: 2⟩ | {I2, I4: 2} |
| I3 | {{I2, I1: 2}, {I2: 2}, {I1: 2}} | ⟨I2: 4, I1: 2⟩, ⟨I1: 2⟩ | {I2, I3: 4}, {I1, I3: 4}, {I2, I1, I3: 2} |
| I1 | {{I2: 4}} | ⟨I2: 4⟩ | {I2, I1: 4} |

# A Special Case: Single Prefix Path in FP-tree

- Suppose a (conditional) FP-tree T has a shared single prefix-path P

- Mining can be decomposed into two parts

  - Reduction of the single prefix path into one node

  - Concatenation of the mining results of the two parts

$$
\begin{array}{c}
\{\} \\
| \\
a_1{:}n_1 \\
| \\
a_2{:}n_2 \\
| \\
a_3{:}n_3 \\
\end{array}
\qquad
\begin{array}{c}
b_1{:}m_1 \quad C_1{:}k_1 \\
\\
C_2{:}k_2 \quad C_3{:}k_3
\end{array}
$$

$\rightarrow$

$$
r_1 =
\begin{array}{c}
\{\} \\
| \\
a_1{:}n_1 \\
| \\
a_2{:}n_2 \\
| \\
a_3{:}n_3
\end{array}
\quad + \quad
\begin{array}{c}
r_1 \\
b_1{:}m_1 \quad C_1{:}k_1 \\
\\
C_2{:}k_2 \quad C_3{:}k_3
\end{array}
$$

# FP-Growth vs. Apriori: Scalability With the Support Threshold



Data set T25I20D10K

# Why Is FP-Growth the Winner?

- Divide-and-conquer:

  - decompose both the mining task and DB according to the frequent patterns obtained so far

  - leads to focused search of smaller databases

- Other factors

  - no candidate generation, no candidate test

  - compressed database: FP-tree structure

  - no repeated scan of entire database

  - basic ops—counting local freq items and building sub FP-tree, no pattern search and matching

# Rule Generation

# 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:

    | | | | |
    |---|---|---|---|
    | 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, | | |
    | ABC $\rightarrow$ D, | ABD $\rightarrow$ C, | ACD $\rightarrow$ B, | BCD $\rightarrow$ A, |

- If |L| = k, then there are $2^k - 2$ candidate association rules (ignoring L $\rightarrow$ $\varnothing$ and $\varnothing$ $\rightarrow$ 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

Low
Confidence
Rule

ABCD=>{ }

BCD=>A   ACD=>B   ABD=>C   ABC=>D

CD=>AB   BD=>AC   BC=>AD   AD=>BC   AC=>BD   AB=>CD

D=>ABC   C=>ABD   B=>ACD   A=>BCD

**Pruned
Rules**

# 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 (number of transactions)
  - Since Apriori makes multiple passes, run time of algorithm may increase with number of transactions
- Average transaction width
  - Transaction width increases with denser data sets
  - This may increase max length of frequent itemsets and traversals of hash tree (number of subsets in a transaction increases with its width)

# **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 | 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 |

- Number of frequent itemsets $= 3 \times \sum\limits_{k=1}^{10} \binom{10}{k}$

- Need a compact representation

# Maximal Frequent Itemset

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

# 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 | 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 |

**Minimum support threshold = 5**

# 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: ?**

# An illustrative example

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

**Transactions**

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

# An illustrative example



**Items**

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

**Transactions**

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

**Support threshold (by count): 4**
**Frequent itemsets: ?**

# An illustrative example

**Items**

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

**Transactions**

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

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

# An illustrative example

**Items**

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

Transactions

**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: ?

# An illustrative example

**Items**

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

**Transactions**

**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



| | 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**

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

**Transactions**

**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

| | A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|---|
| 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**

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

**Transactions**

**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**

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

Transactions

**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.

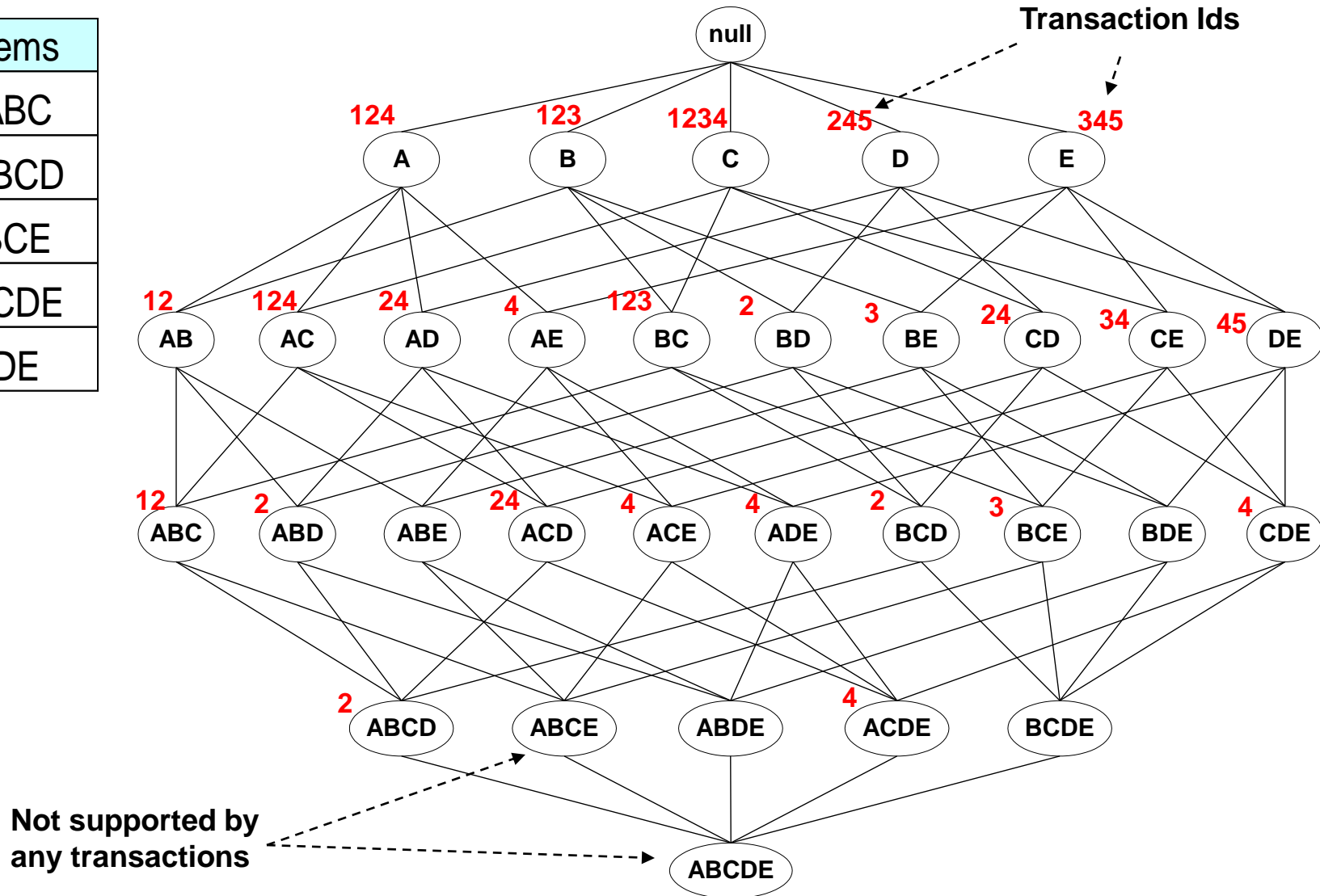- X is not closed if at least one of its immediate supersets has support count as X.

| 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 |

# Maximal vs Closed Itemsets

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



Transaction Ids

Not supported by any transactions

# Maximal vs Closed Frequent Itemsets



**Minimum support = 2**

Closed but not maximal

Closed and maximal

# Frequent = 14

# Closed = 9

# Maximal = 4

# 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 |

# Example 1

**Items**

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

**Transactions**

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

# Example 1

**Items**

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

**Transactions**

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

# Example 2

**Items**

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

**Transactions**

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

# Example 2

**Items**

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

**Transactions**

| Itemsets | Support (counts) | Closed itemsets |
|---|---|---|
| {C} | 3 | ✔ |
| {D} | 2 | |
| {E} | 2 | |
| {C,D} | 2 | |
| {C,E} | 2 | |
| {D,E} | 2 | |
| {C,D,E} | 2 | ✔ |

# Example 3

**Items**

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

**Transactions**

**Closed itemsets:**
**{C,D,E,F}: 2,**
**{C,F}: 3**

# Example 4

**Items**

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

**Transactions**

**Closed itemsets:**
**{C,D,E,F}: 2,**
**{C}: 3,**
**{F}: 3**

# Maximal vs Closed Itemsets



Frequent Itemsets

Closed Frequent Itemsets

Maximal Frequent Itemsets

# 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



a. **What is the number of frequent itemsets for each dataset? Which dataset will produce the most number of frequent itemsets?**

b. **Which dataset will produce the longest frequent itemset?**

c. **Which dataset will produce frequent itemsets with highest maximum support?**

d. **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%)?**

e. **What is the number of maximal frequent itemsets for each dataset? Which dataset will produce the most number of maximal frequent itemsets?**

f. **What is the number of closed frequent itemsets for each dataset? Which dataset will produce the most number of closed 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 <mark>contingency</mark> table

Contingency table

|  | Y | $\overline{Y}$ |  |
|---|---|---|---|
| X | $f_{11}$ | $f_{10}$ | $f_{1+}$ |
| $\overline{X}$ | $f_{01}$ | $f_{00}$ | $f_{o+}$ |
|  | $f_{+1}$ | $f_{+0}$ | N |

$f_{11}$: support of X and Y
$f_{10}$: support of $\underline{X}$ and $\overline{Y}$
$f_{01}$: support of $\underline{X}$ and $\underline{Y}$
$f_{00}$: 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 | $\overline{\text{Coffee}}$ | |
|---|---|---|---|
| Tea | 15 | 5 | 20 |
| $\overline{\text{Tea}}$ | 75 | 5 | 80 |
| | 90 | 10 | 100 |

Association Rule: Tea → Coffee

Confidence ≅ P(Coffee|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 | $\overline{\text{Coffee}}$ | |
|------|------|------|------|
| Tea | 15 | 5 | 20 |
| $\overline{\text{Tea}}$ | 75 | 5 | 80 |
| | 90 | 10 | 100 |

Association Rule: Tea → Coffee

Confidence= P(Coffee|Tea) = 15/20 = 0.75

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

⇒ Note that P(Coffee|$\overline{\text{Tea}}$) = 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

$$\text{confidence}(X \rightarrow Y) = \text{support}(Y)$$

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

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

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

**lift** is used for **rules** while **interest** is used for **itemsets**

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

**Piatesky-Shapiro Measure**

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

**Correlation Analysis**

| | Y | $\overline{Y}$ | |
|---|---|---|---|
| X | $f_{11}$ | $f_{10}$ | $f_{1+}$ |
| $\overline{X}$ | $f_{01}$ | $f_{00}$ | $f_{o+}$ |
| | $f_{+1}$ | $f_{+0}$ | N |

$$\frac{f_{11}f_{00} - f_{01}f_{10}}{\sqrt{f_{1+}f_{+1}f_{0+}f_{+0}}}$$

# Example: Lift/Interest

|      | Coffee | $\overline{\text{Coffee}}$ |      |
|------|--------|--------|------|
| Tea  | 15     | 5      | 20   |
| $\overline{\text{Tea}}$ | 75     | 5      | 80   |
|      | 90     | 10     | 100  |

Association Rule: Tea $\rightarrow$ Coffee

Confidence= P(Coffee|Tea) = 0.75

but P(Coffee) = 0.9

$\Rightarrow$ 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 | $\overline{Y}$ |   |
|---|---|---|---|
| X | 10 | 0 | 10 |
| $\overline{X}$ | 0 | 90 | 90 |
|   | 10 | 90 | 100 |

|   | Y | $\overline{Y}$ |   |
|---|---|---|---|
| X | 90 | 0 | 90 |
| $\overline{X}$ | 0 | 10 | 10 |
|   | 90 | 10 | 100 |

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

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

**Statistical independence:**

**If P(X,Y)=P(X)P(Y)  => Lift = 1**

**There are lots of measures proposed in the literature**

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

95

# 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:

| # | $\phi$ | $\lambda$ | $\alpha$ | $Q$ | $Y$ | $\kappa$ | $M$ | $J$ | $G$ | $s$ | $c$ | $L$ | $V$ | $I$ | $IS$ | $PS$ | $F$ | $AV$ | $S$ | $\zeta$ | $K$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 |

# Property under Variable Permutation

|       | **B** | **B̄** |
|-------|-------|-------|
| **A**  | p     | q     |
| **Ā**  | r     | s     |

$\Longrightarrow$

|       | **A** | **Ā** |
|-------|-------|-------|
| **B**  | p     | r     |
| **B̄**  | 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

# Property under Row/Column Scaling

Grade-Gender Example (Mosteller, 1968):

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

| | 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

Invariant measures:

♦ odds ratio, etc

# Property under Inversion Operation

|   | A | B |   | C | D |   | E | F |
|---|---|---|---|---|---|---|---|---|
| Transaction 1 → | 1 | 0 |   | 0 | 1 |   | 0 | 0 |
| ■ | 0 | 0 |   | 1 | 1 |   | 1 | 0 |
|   | 0 | 0 |   | 1 | 1 |   | 1 | 0 |
| ■ | 0 | 0 |   | 1 | 1 |   | 1 | 0 |
|   | 0 | 1 |   | 1 | 0 |   | 1 | 1 |
| ■ | 0 | 0 |   | 1 | 1 |   | 1 | 0 |
|   | 0 | 0 |   | 1 | 1 |   | 1 | 0 |
| ■ | 0 | 0 |   | 1 | 1 |   | 1 | 0 |
|   | 0 | 0 |   | 1 | 1 |   | 1 | 0 |
| Transaction N → | 1 | 0 |   | 0 | 1 |   | 0 | 0 |
|   | (a) |   |   | (b) |   |   | (c) |   |

# Example: φ-Coefficient

- φ-coefficient is <mark>analogous</mark> to correlation coefficient for continuous variables

| | Y | $\overline{Y}$ | |
|---|---|---|---|
| X | 60 | 10 | 70 |
| $\overline{X}$ | 10 | 20 | 30 |
| | 70 | 30 | 100 |

| | Y | $\overline{Y}$ | |
|---|---|---|---|
| X | 20 | 10 | 30 |
| $\overline{X}$ | 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** | **B̄** |
|---|---|---|
| **A** | p | q |
| **Ā** | r | s |

⟹

| | **B** | **B̄** |
|---|---|---|
| **A** | p | q |
| **Ā** | 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 |
|--------|---------|-----------|---------------|---------|
| $\phi$ | $\phi$-coefficient | Yes | No | No |
| $\alpha$ | odds ratio | Yes | No | Yes |
| $\kappa$ | 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 |
| $\zeta$ | 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

**Support distribution of a retail data set**

# 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

# Cross-Support Patterns



caviar                              milk

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, $X = \{x_1, x_2, \ldots, x_d\}$, with $d$ items, we can define a measure of cross support, named support ratio $r(X)$, for the itemset

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

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

- Given a user-specified threshold $h_c$, an itemset $X$ is a cross support pattern if $r(X) < h_c$.

# Example

| p | q | r |
|---|---|---|
| 0 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 0 | 0 |

- Cross-support patterns can be eliminated
  - {p, q}, {p, r}, and {p, q, r}: $r(X) = 0.2$
- Interesting patterns with low supports are still pruned by support-based pruning strategy
  - {q, r}

- Confidence pruning?
  - {q} → {r}: c = 100%
  - {q} → {p}: c = 80%
  - {r} → {p}: c = 80%

# Confidence and Cross-Support Patterns



caviar

milk

**Observation:**

conf(caviar→milk) is very high

but

conf(milk→caviar) is very low

**Therefore,**

min( conf(caviar→milk),
conf(milk→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 $X = \{x_1, x_2, \ldots, x_d\}$
  - h-confidence is the <u>minimum confidence</u> of any association rule formed from itemset $X$

  - hconf( $X$ ) = min( 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: $X_1 = \{x_1, x_2\}, X_2 = \{x_3, \ldots, x_d\}$

# H-Confidence …

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

$$hconf(X) = \min\left\{\frac{s(X)}{s(x_1)}, \frac{s(X)}{s(x_2)}, \ldots, \frac{s(X)}{s(x_d)}\right\}$$

$$= \frac{s(X)}{\max\{s(x_1), s(x_2), \ldots, s(x_d)\}}$$

# Cross Support and H-confidence

- By the anti-monotone property of support

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

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

$$hconf(X) = \frac{s(X)}{\max\{s(x_1), s(x_2), \ldots, s(x_d)\}}$$

$$\leq \frac{\min\{s(x_1), s(x_2), \ldots, s(x_d)\}}{\max\{s(x_1), s(x_2), \ldots, s(x_d)\}}$$

$$= r(X)$$

Thus, $hconf(X) \leq r(X)$

# Cross Support and H-confidence ...

- Since, $hconf(X) \leq r(X)$, cross support patterns can be eliminated by ensuring that the h-confidence values for the patterns exceed $h_c$, a user set threshold
- Notice that

$$0 \leq hconf(X) \leq r(X) \leq 1$$

- An itemset $X$ is a hyperclique pattern if and only if $hconf(X) > h_c$, where $h_c$ is a given h-confidence threshold
- 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 <u>closed</u> if none of its immediate supersets has the same h-confidence as $X$
  - A hyperclique $X$ is <u>maximal</u> if $\mathrm{hconf}(X) > \mathrm{h_c}$ and none of its immediate supersets, $Y$, have $\mathrm{hconf}(Y) > \mathrm{h_c}$

- Items in a hyperclique cannot have widely different support
  - Allows for more efficient pruning
  - Can be used to find strongly coherent groups of items
    - Words that occur together in documents