# Data Mining
## Association Analysis

马锦华

数据科学与计算机学院

中山大学

# 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

| 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 |

- **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

# Definition: Association Rule

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

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

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

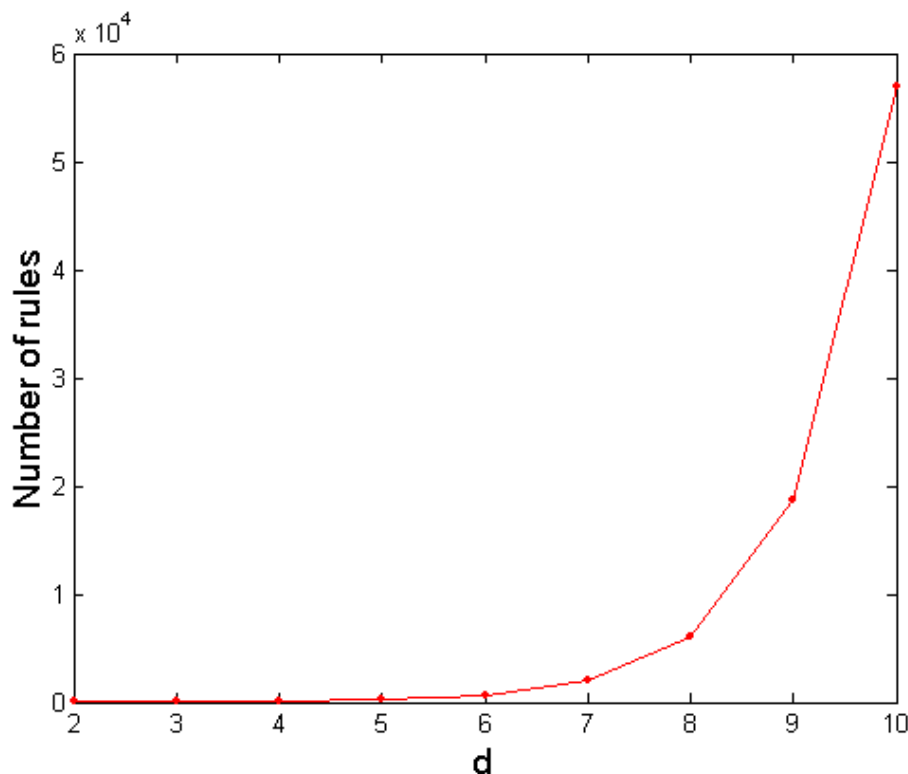$$c = \frac{\sigma(Milk, Diaper, Beer)}{\sigma(Milk, Diaper)} = \frac{2}{3} = 0.67$$

# 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
  - $\Rightarrow$ Computationally prohibitive!

# Computational Complexity

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



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

$$= 3^d - 2^{d+1} + 1$$

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



Given d items, there are $2^d$ possible candidate itemsets

# Frequent Itemset Generation

- Brute-force approach:
  - Each itemset in the lattice is a <span style="color:red">candidate</span> 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) => <span style="color:red">Expensive since M = $2^d$ !!!</span>

# Reducing Number of Candidates

- Apriori principle:
  - If an itemset is frequent, then all of its subsets must also be frequent
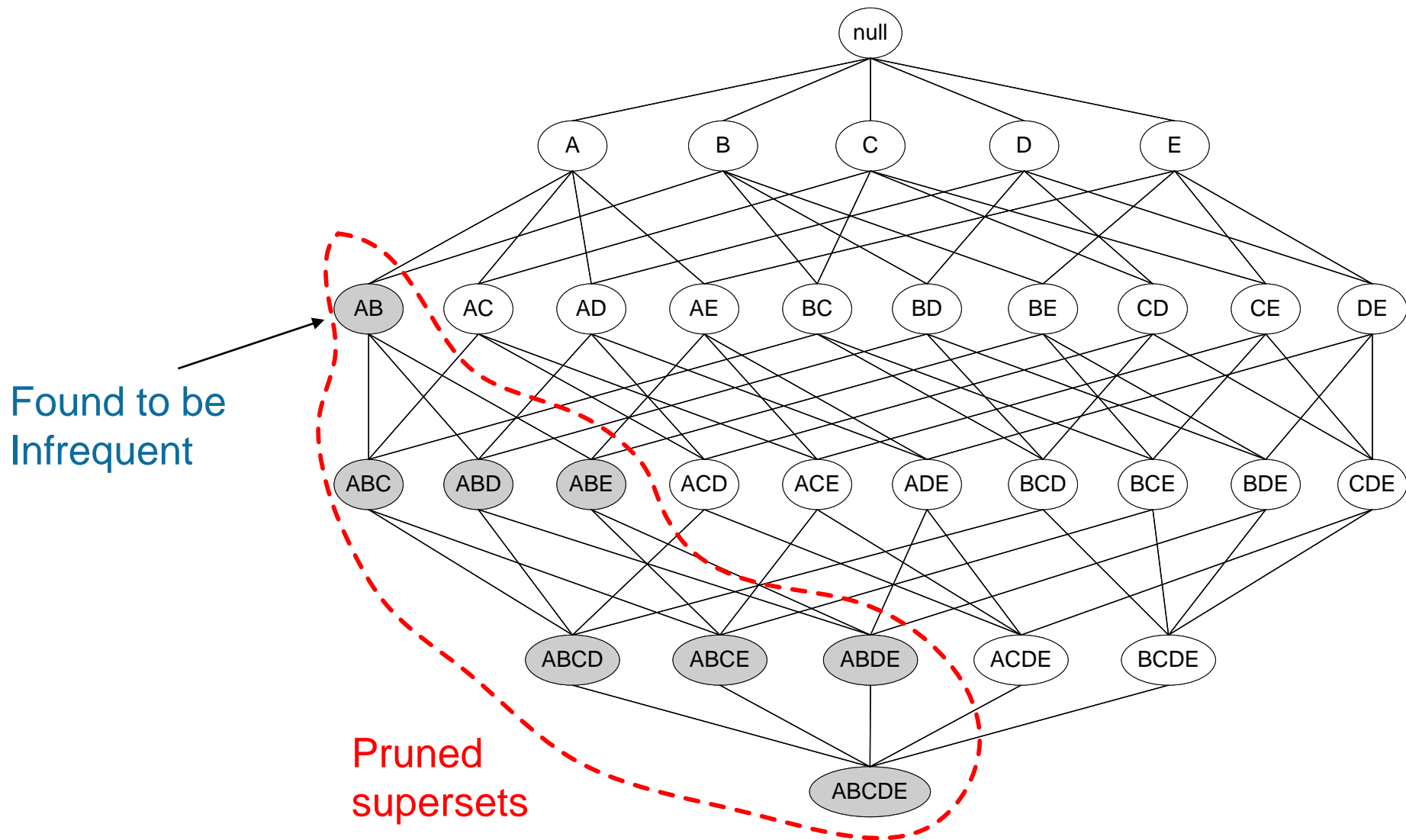
- Apriori principle holds due to the following property of the support measure:

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

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

# Illustrating Apriori Principle

# Illustrating Apriori Principle

| TID | Items |
|-----|-------|
| 1 | Bread, Diaper, 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 | 5 |
| 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

| TID | Items |
|-----|-------|
| 1 | Bread, Diaper, 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 | 5 |
| 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} | 4 |
| {Beer,Milk} | 2 |
| {Diaper,Milk} | 4 |
| {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}** | **4** |
| **{Milk,Beer}** | **2** |
| **{Milk,Diaper}** | **4** |
| **{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,
$$6 + 6 + 4 = 16$$

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

Triplets (3-itemsets)

# 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)

Triplets (3-itemsets)

| Itemset | Count |
|---------|-------|
| { Beer, Diaper, Milk} | 2 |
| { Beer,Bread, Diaper} | 2 |
| **{Bread, Diaper, Milk}** | **3** |
| {Beer, Bread, Milk} | 1 |

Minimum Support = 3

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

# Apriori Algorithm

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

- Algorithm
  - Let k=1
  - Generate $F_1$ = {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 => $F_{k+1}$

# 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 $L \rightarrow \varnothing$ and $\varnothing \rightarrow L$)
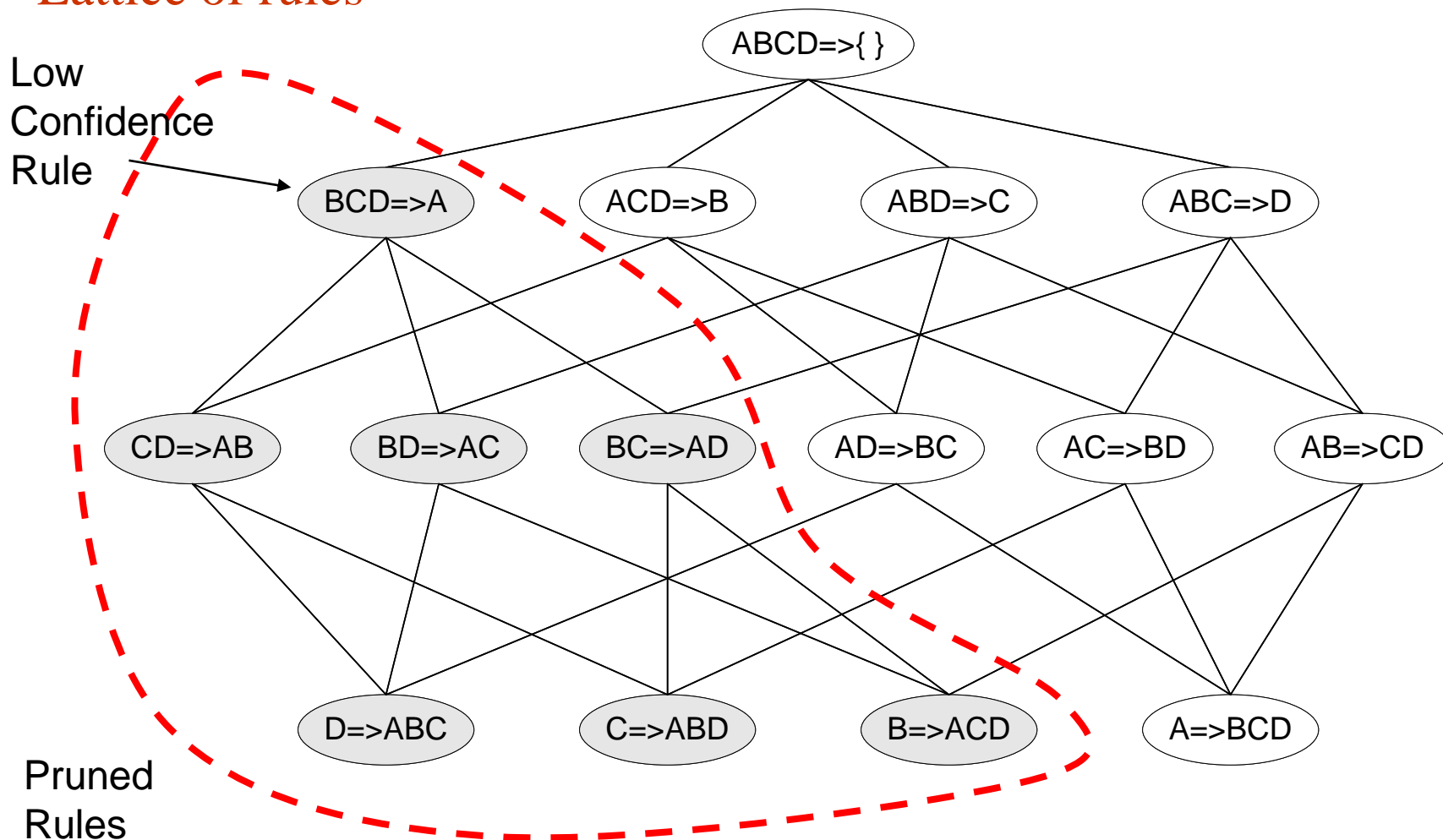
# Rule Generation

- 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 right hand side of the rule

# Rule Generation for Apriori Algorithm
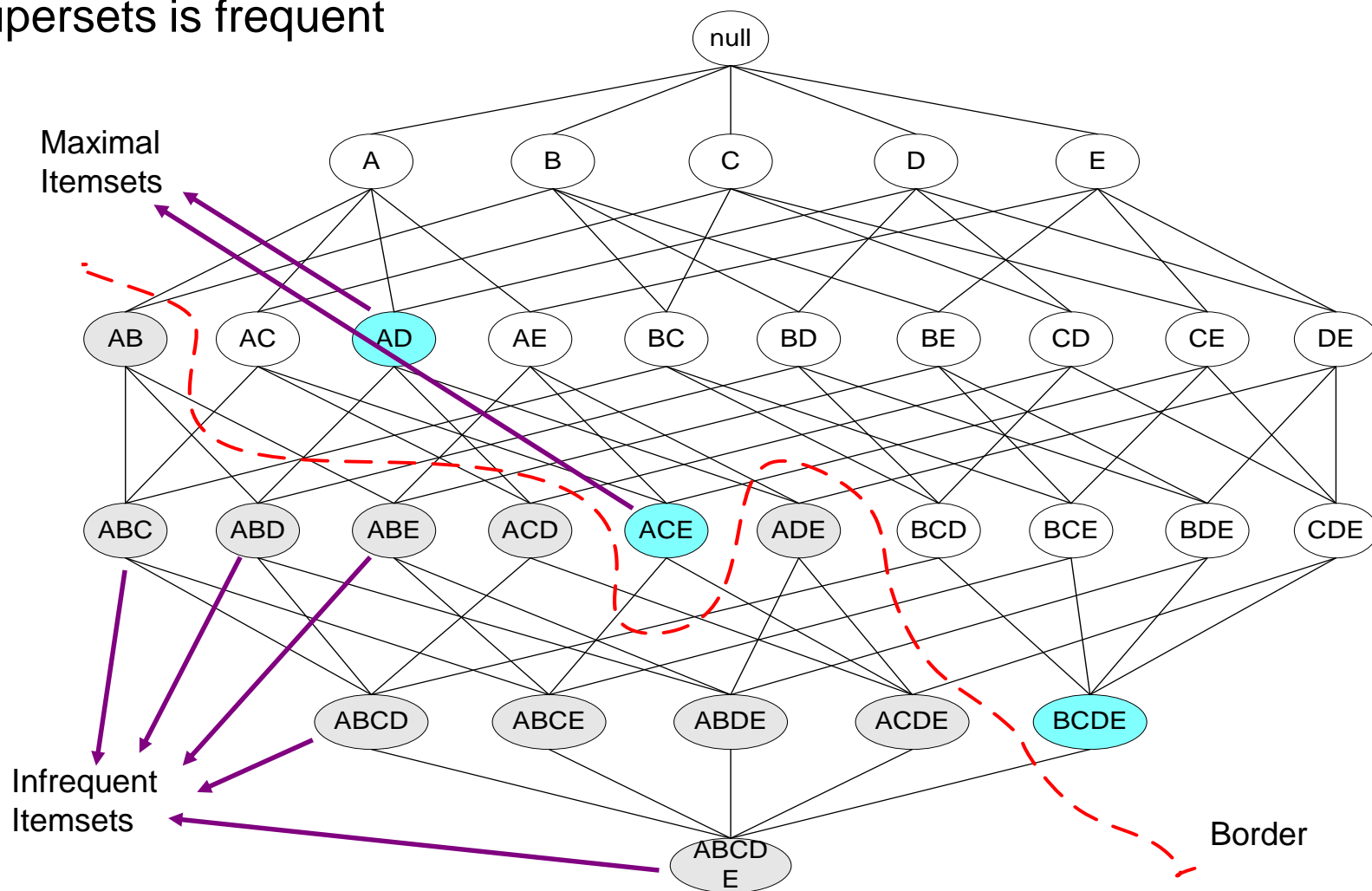
Lattice of rules

# 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 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)

# Maximal Frequent Itemset

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

# 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}

# 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:
All subsets of {C,D,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}
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

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

Items (column header)

Transactions (row label)

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

| | 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: ?

# An illustrative example



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

null

124 — A
123 — B
1234 — C
245 — D
345 — E

12 — AB
124 — AC
24 — AD
4 — AE
123 — BC
2 — BD
3 — BE
24 — CD
34 — CE
45 — DE

12 — ABC
2 — ABD
ABE
24 — ACD
4 — ACE
4 — ADE
2 — BCD
3 — BCE
BDE
4 — CDE

2 — ABCD
ABCE
ABDE
4 — ACDE
BCDE

ABCDE

# Closed = 9

# Maximal = 4

# 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}, {C,F}

# 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}, {C}, {F}

# Maximal vs Closed 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
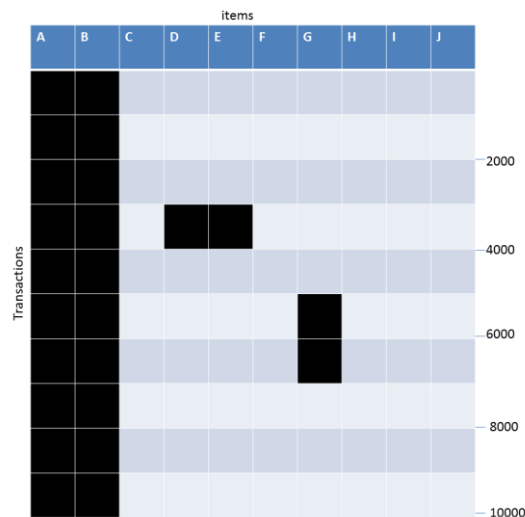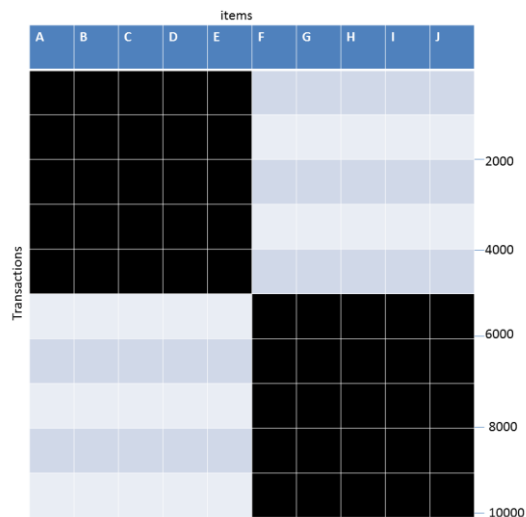


a. What is the number of frequent itemsets for each dataset? Which dataset will produce the most number of frequent itemsets?

# 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
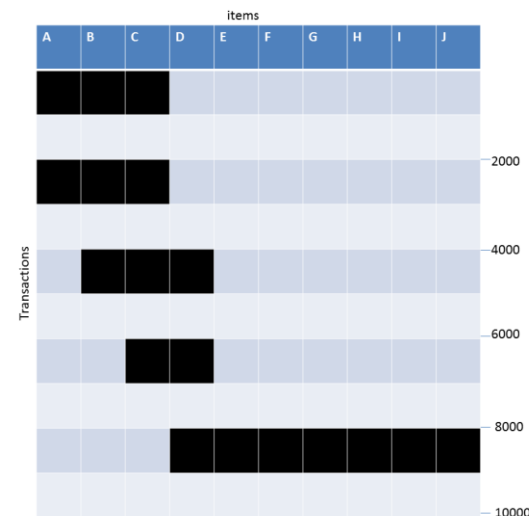


b. Which dataset will produce the longest frequent itemset?
c. Which dataset will produce frequent itemsets with highest maximum support?

# 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
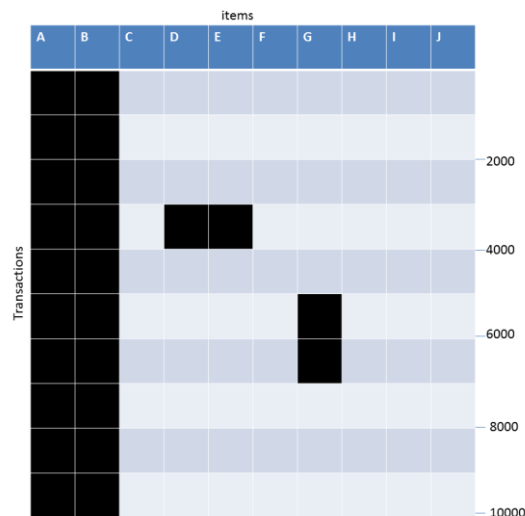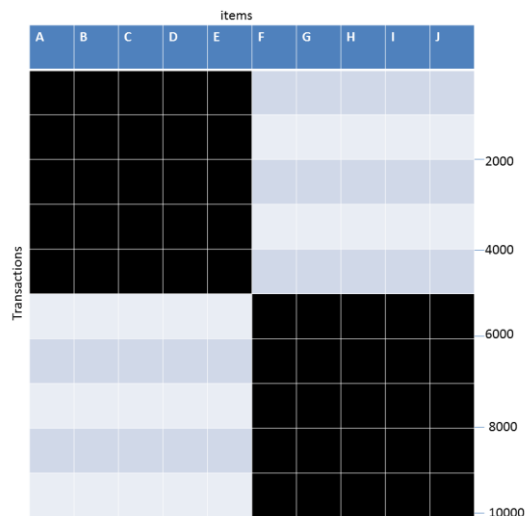


d. Which dataset will produce frequent itemsets containing items with widely varying support levels (i.e., itemsets containing items with mixed support, ranging from 20% to more than 70%)?

# 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



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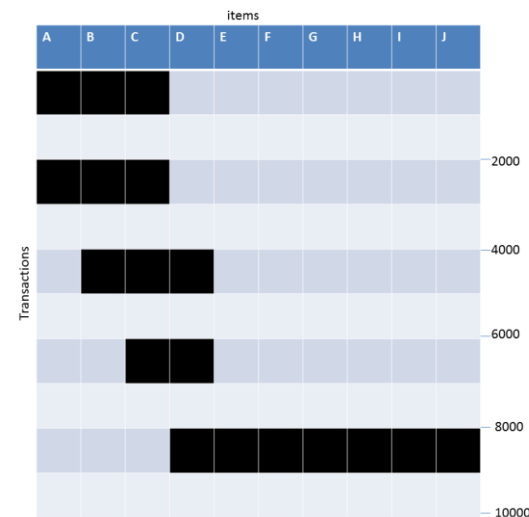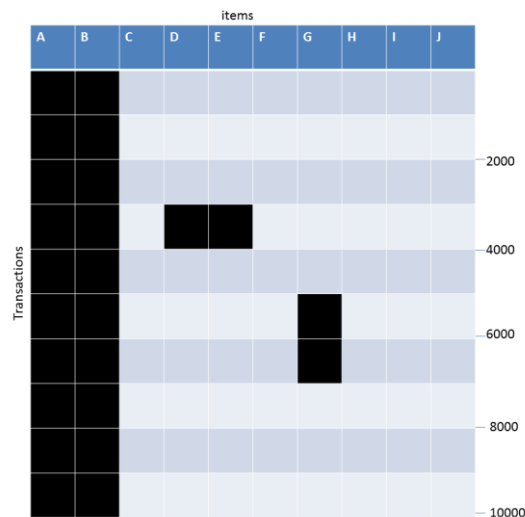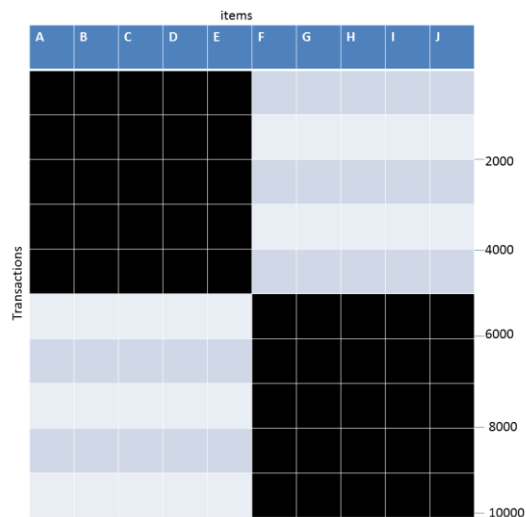


e. What is the number of closed frequent itemsets for each dataset? Which dataset will produce the most number of closed frequent itemsets?

# 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 contingency table

<span style="color:red">Contingency table</span>

|  | 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 $\overline{Y}$

<span style="color:red">Used to define various measures</span>

◆ support, confidence, Gini, entropy, etc.

# Drawback of Confidence

| Customers | 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 $\cong$ 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 $\rightarrow$ 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!

$\Rightarrow$ 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

  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

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

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

lift is used for rules while interest is used for itemsets

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

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

# Example: Lift/Interest

|      | 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) = 0.75

but P(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 | $\overline{Y}$ | |
|---|---|---|---|
| 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 $Q$ | $\dfrac{P(A,B)P(\overline{AB})-P(A,\overline{B})P(\overline{A},B)}{P(A,B)P(\overline{AB})+P(A,\overline{B})P(\overline{A},B)}=\dfrac{\alpha-1}{\alpha+1}$ |
| 5 | Yule's $Y$ | $\dfrac{\sqrt{P(A,B)P(\overline{AB})}-\sqrt{P(A,\overline{B})P(\overline{A},B)}}{\sqrt{P(A,B)P(\overline{AB})}+\sqrt{P(A,\overline{B})P(\overline{A},B)}}=\dfrac{\sqrt{\alpha}-1}{\sqrt{\alpha}+1}$ |
| 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 $(J)$ | $\max\left(P(A,B)\log(\frac{P(B|A)}{P(B)})+P(A\overline{B})\log(\frac{P(\overline{B}|A)}{P(\overline{B})}),\right.$ $\left.P(A,B)\log(\frac{P(A|B)}{P(A)})+P(\overline{A}B)\log(\frac{P(\overline{A}|B)}{P(A)})\right)$ |
| 9 | Gini index $(G)$ | $\max\left(P(A)[P(B|A)^2+P(\overline{B}|A)^2]+P(\overline{A})[P(B|\overline{A})^2+P(\overline{B}|\overline{A})^2]\right.$ $-P(B)^2-P(\overline{B})^2,$ $P(B)[P(A|B)^2+P(\overline{A}|B)^2]+P(\overline{B})[P(A|\overline{B})^2+P(\overline{A}|\overline{B})^2]$ $\left.-P(A)^2-P(\overline{A})^2\right)$ |
| 10 | Support $(s)$ | $P(A,B)$ |
| 11 | Confidence $(c)$ | $\max(P(B|A),P(A|B))$ |
| 12 | Laplace $(L)$ | $\max\left(\frac{NP(A,B)+1}{NP(A)+2},\frac{NP(A,B)+1}{NP(B)+2}\right)$ |
| 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 $(F)$ | $\max\left(\frac{P(B|A)-P(B)}{1-P(B)},\frac{P(A|B)-P(A)}{1-P(A)}\right)$ |
| 18 | Added Value $(AV)$ | $\max(P(B|A)-P(B),P(A|B)-P(A))$ |
| 19 | Collective strength $(S)$ | $\dfrac{P(A,B)+P(\overline{AB})}{P(A)P(B)+P(\overline{A})P(\overline{B})}\times\dfrac{1-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A,B)-P(\overline{AB})}$ |
| 20 | Jaccard $(\zeta)$ | $\dfrac{P(A,B)}{P(A)+P(B)-P(A,B)}$ |
| 21 | Klosgen $(K)$ | $\sqrt{P(A,B)}\max(P(B|A)-P(B),P(A|B)-P(A))$ |

# 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):

|  | 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

| | 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 analogous 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 |

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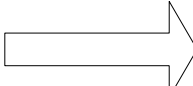$$= 0.5238$$

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

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 |

# References

- P.-N. Tan, M. Steinbach, V. Kumar: Introduction to data mining, Second Edition, https://www-users.cs.umn.edu/~kumar001/dmbook/index.php