# CSE 5243 INTRO. TO DATA MINING

## Mining Frequent Patterns and Associations: Basic Concepts

### Yu Su, CSE@The Ohio State University

# Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

☐ Basic Concepts

☐ Efficient Pattern Mining Methods

☐ Pattern Evaluation

☐ Summary

# Pattern Discovery: Basic Concepts

- What Is Pattern Discovery?   Why Is It Important?

- Basic Concepts: Frequent Patterns and Association Rules

- Compressed Representation: Closed Patterns and Max-Patterns

# What Is Pattern Discovery?

- Motivating examples:
  - What products were often purchased together?
  - What are the subsequent purchases after buying an iPad?
  - What code segments likely contain copy-and-paste bugs?
  - What word sequences likely form phrases in this corpus?

# What Is Pattern Discovery?

- Motivation examples:

  - What products were often purchased together?

  - What are the subsequent purchases after buying an iPad?

  - What code segments likely contain copy-and-paste bugs?

  - What word sequences likely form phrases in this corpus?

- What are patterns?

  - Patterns: A set of items, subsequences, or substructures that occur frequently together (or strongly correlated) in a data set

  - Patterns represent intrinsic and important properties of datasets

# What Is Pattern Discovery?

- Motivation examples:

  - What products were often purchased together?

  - What are the subsequent purchases after buying an iPad?

  - What code segments likely contain copy-and-paste bugs?

  - What word sequences likely form phrases in this corpus?

- What are patterns?

  - Patterns: A set of items, subsequences, or substructures that occur frequently together (or strongly correlated) in a data set

  - Patterns represent intrinsic and important properties of datasets

- Pattern discovery: Uncovering patterns from massive data sets

# Pattern Discovery: Why Is It Important?

- Finding inherent regularities in a data set

- Foundation for many essential data mining tasks

  - Association, correlation, and causality analysis

  - Mining sequential, structural (e.g., sub-graph) patterns

  - Pattern analysis in spatiotemporal, multimedia, time-series, and stream data

  - Classification: Discriminative pattern-based analysis

  - Cluster analysis: Pattern-based subspace clustering

# Pattern Discovery: Why Is It Important?

- Finding inherent regularities in a data set
- Foundation for many essential data mining tasks
  - Association, correlation, and causality analysis
  - Mining sequential, structural (e.g., sub-graph) patterns
  - Pattern analysis in spatiotemporal, multimedia, time-series, and stream data
  - Classification: Discriminative pattern-based analysis
  - Cluster analysis: Pattern-based subspace clustering
- Broad applications
  - Market basket analysis, cross-marketing, catalog design, sale campaign analysis, Web log analysis, biological sequence analysis

# Basic Concepts: k-Itemsets and Their Supports

- Itemset: A set of one or more items

# Basic Concepts: k-Itemsets and Their Supports

- Itemset: A set of one or more items

- k-itemset:  X = {$x_1$, ..., $x_k$}
  - Ex. {Beer, Nuts, Diaper} is a 3-itemset

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

# Basic Concepts: k-Itemsets and Their Supports

- Itemset: A set of one or more items

- k-itemset:  X = $\{x_1, ..., x_k\}$
  - Ex. {Beer, Nuts, Diaper} is a 3-itemset

- (*absolute*) *support* (*count*) of X, sup{X}: Frequency or the number of occurrences of an itemset X

  - Ex.  sup{Beer} = 3
  - Ex.  sup{Diaper} = 4
  - Ex.  sup{Beer, Diaper} = 3
  - Ex.  sup{Beer, Eggs} = 1

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

# Basic Concepts: k-Itemsets and Their Supports

- Itemset: A set of one or more items

- k-itemset:  X = $\{x_1, ..., x_k\}$
  - Ex. {Beer, Nuts, Diaper} is a 3-itemset

- (*absolute*) *support* (*count*) of X, sup{X}: Frequency or the number of occurrences of an itemset X
  - Ex.  sup{Beer} = 3
  - Ex.  sup{Diaper} = 4
  - Ex.  sup{Beer, Diaper} = 3
  - Ex.  sup{Beer, Eggs} = 1

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

- (*relative*) *support*, *s*{X}:  The fraction of transactions that contains X (i.e., the probability that a transaction contains X)
  - Ex.  s{Beer} = 3/5 = 60%
  - Ex.  s{Diaper} = 4/5 = 80%
  - Ex.  s{Beer, Eggs} = 1/5 = 20%

# Basic Concepts: Frequent Itemsets (Patterns)

- An itemset (or a pattern) X is *frequent* if the support of X is no less than a *minsup* threshold $\sigma$

# Basic Concepts: Frequent Itemsets (Patterns)

☐ An itemset (or a pattern) X is *frequent* if the support of X is no less than a *minsup* threshold σ

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

☐ Let σ = *50%*  (σ: *minsup* threshold)
For the given 5-transaction dataset
- ☐ All the frequent 1-itemsets:
  - ■ Beer: 3/5 (60%); Nuts: 3/5 (60%)
  - ■ Diaper: 4/5 (80%); Eggs: 3/5 (60%)
- ☐ All the frequent 2-itemsets:
  - ■ {Beer, Diaper}: 3/5 (60%)
- ☐ All the frequent 3-itemsets?
  - ■ None

14

# Basic Concepts: Frequent Itemsets (Patterns)

An itemset (or a pattern) X is *frequent* if the support of X is no less than a *minsup* threshold σ

Let σ = *50%*  (σ: *minsup* threshold)
For the given 5-transaction dataset
- All the frequent 1-itemsets:
  - Beer: 3/5 (60%); Nuts: 3/5 (60%)
  - Diaper: 4/5 (80%); Eggs: 3/5 (60%)
- All the frequent 2-itemsets:
  - {Beer, Diaper}: 3/5 (60%)
- All the frequent 3-itemsets?
  - None

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

- Do these itemsets (shown on the left) form the complete set of frequent *k*-itemsets (patterns) for any *k*?
- **Observation**:  We may need an efficient method to mine a complete set of frequent patterns

# From Frequent Itemsets to Association Rules

- Comparing with itemsets, rules can be more telling
  - Ex. *Diaper* → *Beer*
    - *Buying diapers may likely lead to buying beers*

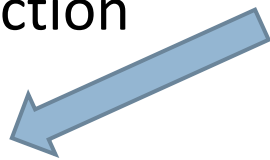# From Frequent Itemsets to Association Rules

- □ Ex. *Diaper → Beer: Buying diapers may likely lead to buying beers*

□ How strong is this rule? (support, confidence)

- □ Measuring association rules: $X → Y$ (s, c)
  - Both $X$ and $Y$ are itemsets

# From Frequent Itemsets to Association Rules

□ Ex. *Diaper → Beer: Buying diapers may likely lead to buying beers*

□ How strong is this rule?  (support, confidence)

■ Measuring association rules:  $X → Y$ (s, c)

■ Both $X$ and $Y$ are itemsets

■ Support, $s$: The probability that a transaction contains $X \cup Y$

■ Ex. s{Diaper, Beer} = 3/5 = 0.6 (i.e., 60%)

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

# From Frequent Itemsets to Association Rules

- Ex. *Diaper → Beer : Buying diapers may likely lead to buying beers*

☐ How strong is this rule?  (support, confidence)

- Measuring association rules:  $X → Y$ (s, c)
  - Both *X* and *Y* are itemsets

- Support, *s*: The probability that a transaction contains $X \cup Y$
  - Ex. s{Diaper, Beer} = 3/5 = 0.6 (i.e., 60%)

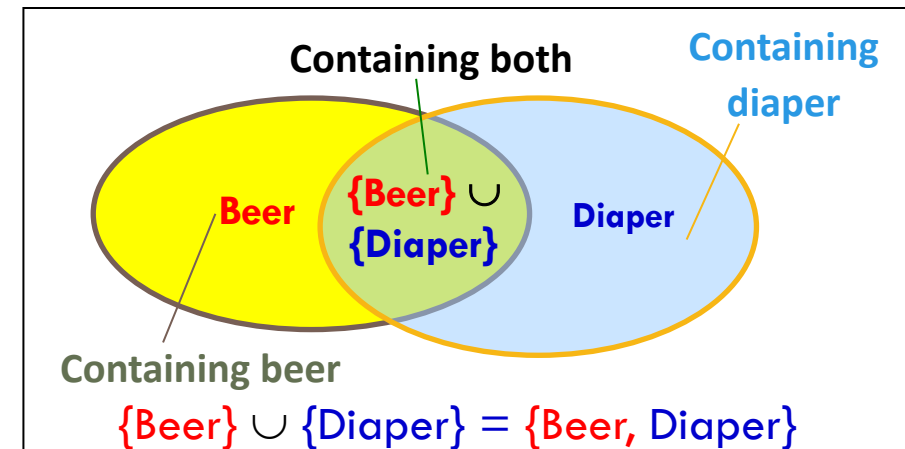- Confidence, *c: The conditional probability* that a transaction containing X also contains *Y*
  - Calculation: $c = sup(X \cup Y) / sup(X)$
  - Ex. $c = sup\{Diaper, Beer\}/sup\{Diaper\} = ¾ = 0.75$

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



Containing both

Containing diaper

Beer

{Beer} ∪ {Diaper}

Diaper

Containing beer

{Beer} ∪ {Diaper} = {Beer, Diaper}

# Mining Frequent Itemsets and Association Rules

- **Association rule mining**
  - Given two thresholds: *minsup, minconf*
  - Find <span style="color:red">all</span> of the rules, $X \rightarrow Y$ (s, c)
    - such that, s ≥ *minsup* and  c ≥ *minconf*

# Mining Frequent Itemsets and Association Rules

**Association rule mining**

- Given two thresholds: *minsup, minconf*
- Find <span style="color:red">all</span> of the rules, $X \rightarrow Y$ (s, c)
  - such that, s ≥ *minsup* and c ≥ *minconf*

- Let *minsup = 50%*
  - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
  - Freq. 2-itemsets: {Beer, Diaper}: 3

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

# Mining Frequent Itemsets and Association Rules

- **Association rule mining**
  - Given two thresholds: *minsup, minconf*
  - Find **all** of the rules, $X \rightarrow Y$ (s, c)
    - such that, s ≥ *minsup* and  c ≥ *minconf*

- Let  *minsup* = 50%
  - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
  - Freq. 2-itemsets:  {Beer, Diaper}: 3

- Let *minconf* = 50%
  - *Beer $\rightarrow$ Diaper  (60%, 100%)*
  - *Diaper $\rightarrow$ Beer  (60%, 75%)*

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

# Mining Frequent Itemsets and Association Rules

**Association rule mining**

- Given two thresholds: *minsup, minconf*
- Find **all** of the rules, *X* → *Y* (s, c)
  - such that, s ≥ *minsup* and c ≥ *minconf*

- Let *minsup* = 50%
  - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
  - Freq. 2-itemsets: {Beer, Diaper}: 3

- Let *minconf* = 50%
  - *Beer → Diaper* (60%, 100%)
  - *Diaper → Beer* (60%, 75%)

(Q: Are these all rules?)

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

# Mining Frequent Itemsets and Association Rules

□ **Association rule mining**

    ▫ Given two thresholds: *minsup, minconf*

    ▫ Find **all** of the rules, $X \rightarrow Y$ (s, c)

        ■ such that, s ≥ *minsup* and  c ≥ *minconf*

□   Let  *minsup* = 50%

    ❑   Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3

    ❑   Freq. 2-itemsets:  {Beer, Diaper}: 3

□   Let *minconf* = 50%

    ❑   *Beer → Diaper*  (60%, 100%)

    ❑   *Diaper → Beer*  (60%, 75%)

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

□  **Observations:**

    ❑   Mining association rules and mining frequent patterns are very close problems

    ❑   Scalable methods are needed for mining large datasets

# Association Rule Mining: two-step process

In general, association rule mining can be viewed as a two-step process:

1. **Find all frequent itemsets:** By definition, each of these itemsets will occur at least as frequently as a predetermined minimum support count, $min\_sup$.

2. **Generate strong association rules from the frequent itemsets:** By definition, these rules must satisfy minimum support and minimum confidence.

Because the second step is much less costly than the first, the overall performance of mining association rules is determined by the first step.

# Generating Association Rules from Frequent Patterns

☐ Recall that:

$$confidence(A \Rightarrow B) = P(B|A) = \frac{support\_count(A \cup B)}{support\_count(A)}.$$

☐ Once we mined frequent patterns, association rules can be generated as follows:

- For each frequent itemset $l$, generate all nonempty subsets of $l$.

- For every nonempty subset $s$ of $l$, output the rule "$s \Rightarrow (l - s)$" if $\frac{support\_count(l)}{support\_count(s)} \geq min\_conf$, where $min\_conf$ is the minimum confidence threshold.

# Generating Association Rules from Frequent Patterns

□ Recall that:

$$confidence(A \Rightarrow B) = P(B|A) = \frac{support\_count(A \cup B)}{support\_count(A)}$$

□ Once we mined frequent patterns, association rules can be generated as follows:

- For each frequent itemset $l$, generate all nonempty subsets of $l$.

- For every nonempty subset $s$ of $l$, output the rule "$s \Rightarrow (l - s)$" if $\frac{support\_count(l)}{support\_count(s)} \geq min\_conf$, where $min\_conf$ is the minimum confidence threshold.

Because $l$ is a frequent itemset, each rule automatically satisfies the minimum support requirement.

# Example: Generating Association Rules

**Generating association rules.** Let's try an example based on the transactional data for *AllElectronics* shown in Table 6.1. The data contain frequent itemset $X = \{I1, I2, I5\}$. What are the association rules that can be generated from $X$? The nonempty subsets of $X$ are $\{I1, I2\}$, $\{I1, I5\}$, $\{I2, I5\}$, $\{I1\}$, $\{I2\}$, and $\{I5\}$. The resulting association rules are as shown below, each listed with its confidence:

$$\{I1, I2\} \Rightarrow I5, \qquad confidence = 2/4 = 50\%$$
$$\{I1, I5\} \Rightarrow I2, \qquad confidence = 2/2 = 100\%$$
$$\{I2, I5\} \Rightarrow I1, \qquad confidence = 2/2 = 100\%$$
$$I1 \Rightarrow \{I2, I5\}, \qquad confidence = 2/6 = 33\%$$
$$I2 \Rightarrow \{I1, I5\}, \qquad confidence = 2/7 = 29\%$$
$$I5 \Rightarrow \{I1, I2\}, \qquad confidence = 2/2 = 100\%$$

If minimum confidence threshold: 70%, what will be output?

# Challenge: There Are Too Many Frequent Patterns!

- A long pattern contains a combinatorial number of sub-patterns

- How many frequent itemsets does the following $TDB_1$ contain?

  - $TDB_1$:     $T_1$: $\{a_1, \ldots, a_{50}\}$;  $T_2$: $\{a_1, \ldots, a_{100}\}$

  - Assuming (absolute) *minsup* = 1

  - Let's give it a try…

  1-itemsets:  $\{a_1\}$: 2, $\{a_2\}$: 2, …, $\{a_{50}\}$: 2, $\{a_{51}\}$: 1, …, $\{a_{100}\}$: 1,

  2-itemsets: $\{a_1, a_2\}$: 2, …, $\{a_1, a_{50}\}$: 2, $\{a_1, a_{51}\}$: 1 …, …, $\{a_{99}, a_{100}\}$: 1,

  …, …, …, …

  99-itemsets: $\{a_1, a_2, \ldots, a_{99}\}$: 1, …, $\{a_2, a_3, \ldots, a_{100}\}$: 1

  100-itemset: $\{a_1, a_2, \ldots, a_{100}\}$: 1

# Challenge: There Are Too Many Frequent Patterns!

- A long pattern contains a combinatorial number of sub-patterns

- How many frequent itemsets does the following $TDB_1$ contain?

  - $TDB_1$:      $T_1$: $\{a_1, ..., a_{50}\}$;   $T_2$: $\{a_1, ..., a_{100}\}$
  - Assuming (absolute) *minsup* = 1
  - Let's give it a try…

  1-itemsets: $\{a_1\}$: 2, $\{a_2\}$: 2, …, $\{a_{50}\}$: 2, $\{a_{51}\}$: 1, …, $\{a_{100}\}$: 1,

  2-itemsets: $\{a_1, a_2\}$: 2, …, $\{a_1, a_{50}\}$: 2, $\{a_1, a_{51}\}$: 1 …, …, $\{a_{99}, a_{100}\}$: 1,

  …, …, …, …

  99-itemsets: $\{a_1, a_2, ..., a_{99}\}$: 1, …, $\{a_2, a_3, ..., a_{100}\}$: 1

  100-itemset: $\{a_1, a_2, ..., a_{100}\}$: 1

- The total number of frequent itemsets:

$$\binom{100}{1} + \binom{100}{2} + \binom{100}{3} + \cdots + \binom{100}{100} = 2^{100} - 1$$

Too huge a set for any one to compute or store!

# Expressing Patterns in Compressed Form: Closed Patterns

- How to handle such a challenge?

- Solution 1: **Closed patterns:**  A pattern (itemset) X is <span style="color:red">closed</span> if X is *frequent,* and there exists *no super-pattern* Y ⊃ X, *with the same support* as X

# Expressing Patterns in Compressed Form: Closed Patterns

☐ How to handle such a challenge?

☐ Solution 1: **Closed patterns**:  A pattern (itemset) X is <span style="color:red">closed</span> if X is *frequent,* and there exists *no super-pattern* Y ⊃ X, *with the same support* as X

 ☐ Let Transaction DB $TDB_1$:   $T_1$: $\{a_1, …, a_{50}\}$;  $T_2$: $\{a_1, …, a_{100}\}$

 ☐ Suppose *minsup* = 1. How many closed patterns does $TDB_1$ contain?

   ■ Two:  $P_1$: "$\{a_1, …, a_{50}\}$: 2";  $P_2$: "$\{a_1, …, a_{100}\}$: 1"

Why?

# Expressing Patterns in Compressed Form: Closed Patterns

- How to handle such a challenge?

- Solution 1: **Closed patterns**: A pattern (itemset) X is closed if X is *frequent,* and there exists *no super-pattern* $Y \supset X$, *with the same support* as X

  - Let Transaction DB $TDB_1$: $T_1$: $\{a_1, ..., a_{50}\}$; $T_2$: $\{a_1, ..., a_{100}\}$

  - Suppose *minsup* = 1. How many closed patterns does $TDB_1$ contain?

    - Two: $P_1$: "$\{a_1, ..., a_{50}\}$: 2"; $P_2$: "$\{a_1, ..., a_{100}\}$: 1"

- Closed pattern is a lossless compression of frequent patterns

  - Reduces the # of patterns but does not lose the support information!

  - You will still be able to say: "$\{a_2, ..., a_{40}\}$: 2", "$\{a_5, a_{51}\}$: 1"

# Expressing Patterns in Compressed Form: Max-Patterns

- Solution 2: **Max-patterns**: A pattern X is a max-pattern if X is frequent and there exists no frequent super-pattern Y ⊃ X

# Expressing Patterns in Compressed Form: Max-Patterns

- Solution 2: **Max-patterns**: A pattern X is a max-pattern if X is frequent and there exists no frequent super-pattern $Y \supset X$

- Difference with closed-patterns?

  - Do not care about the real support of the sub-patterns of a max-pattern

  - Let Transaction DB $TDB_1$: $T_1$: $\{a_1, \ldots, a_{50}\}$; $T_2$: $\{a_1, \ldots, a_{100}\}$

  - Suppose *minsup* = 1. How many max-patterns does $TDB_1$ contain?

    - One: P: "$\{a_1, \ldots, a_{100}\}$: 1"

Why?

# Expressing Patterns in Compressed Form: Max-Patterns

- Solution 2: **Max-patterns**: A pattern X is a max-pattern if X is frequent and there exists no frequent super-pattern $Y \supset X$

- Difference with close-patterns?

  - Do not care about the real support of the sub-patterns of a max-pattern

  - Let Transaction DB $TDB_1$: $T_1$: $\{a_1, \ldots, a_{50}\}$; $T_2$: $\{a_1, \ldots, a_{100}\}$

  - Suppose *minsup* = 1. How many max-patterns does $TDB_1$ contain?

    - One: P: "$\{a_1, \ldots, a_{100}\}$: 1"

- Max-pattern is a lossy compression!

  - We only know $\{a_1, \ldots, a_{40}\}$ is frequent

  - But we do not know the real support of $\{a_1, \ldots, a_{40}\}$, …, any more!

  - Thus in many applications, closed-patterns are more desirable than max-patterns

# Example

**Closed and maximal frequent itemsets.** Suppose that a transaction database has only two transactions: $\{\langle a_1, a_2, \ldots, a_{100} \rangle; \langle a_1, a_2, \ldots, a_{50} \rangle\}$. Let the minimum support count threshold be $min\_sup = 1$. We find two closed frequent itemsets and their support counts, that is, $\mathcal{C} = \{\{a_1, a_2, \ldots, a_{100}\} : 1; \{a_1, a_2, \ldots, a_{50}\} : 2\}$. There is only one maximal frequent itemset: $\mathcal{M} = \{\{a_1, a_2, \ldots, a_{100}\} : 1\}$. Notice that we cannot include $\{a_1, a_2, \ldots, a_{50}\}$ as a maximal frequent itemset because it has a frequent super-set, $\{a_1, a_2, \ldots, a_{100}\}$. Compare this to the above, where we determined that there are $2^{100} - 1$ frequent itemsets, which is too huge a set to be enumerated!

{all frequent patterns} $\supseteq$ {closed frequent patterns} $\supseteq$ {max frequent patterns}

# Example

**Closed and maximal frequent itemsets.** Suppose that a transaction database has only two transactions: $\{\langle a_1, a_2, \ldots, a_{100}\rangle; \langle a_1, a_2, \ldots, a_{50}\rangle\}$. Let the minimum support count threshold be $min\_sup = 1$. We find two closed frequent itemsets and their support counts, that is, $\mathcal{C} = \{\{a_1, a_2, \ldots, a_{100}\} : 1; \{a_1, a_2, \ldots, a_{50}\} : 2\}$. There is only one maximal frequent itemset: $\mathcal{M} = \{\{a_1, a_2, \ldots, a_{100}\} : 1\}$. Notice that we cannot include $\{a_1, a_2, \ldots, a_{50}\}$ as a maximal frequent itemset because it has a frequent super-set, $\{a_1, a_2, \ldots, a_{100}\}$. Compare this to the above, where we determined that there are $2^{100} - 1$ frequent itemsets, which is too huge a set to be enumerated!

The set of closed-patterns contains complete information regarding the frequent itemsets.

# Quiz

□ Given closed frequent itemsets:

$$C = \{ \{a1, a2, …, a100\}: 1; \quad \{a1, a2, …, a50\}: 2 \}$$

Is {a2, a45} frequent? Can we know its support?

# Quiz (Cont'd)

- Given maximal frequent itemset:

$$M = \{\{a1, a2, \ldots, a100\}: 1\}$$

What is the support of {a8, a55}?

# Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

☐ Basic Concepts

☐ Efficient Pattern Mining Methods

  ☐ The Apriori Algorithm

  ☐ Application in Classification

☐ Pattern Evaluation

☐ Summary

# Efficient Pattern Mining Methods

- The Downward Closure Property of Frequent Patterns

- <span style="color:red">The Apriori Algorithm</span>

- Extensions or Improvements of Apriori

- Mining Frequent Patterns by Exploring Vertical Data Format

- FPGrowth:  A Frequent Pattern-Growth Approach

- Mining Closed Patterns

# The Downward Closure Property of Frequent Patterns

- Observation: From $TDB_1$: $T_1$: $\{a_1, \ldots, a_{50}\}$; $T_2$: $\{a_1, \ldots, a_{100}\}$
  - We get a frequent itemset: $\{a_1, \ldots, a_{50}\}$
  - Also, its subsets are all frequent: $\{a_1\}$, $\{a_2\}$, $\ldots$, $\{a_{50}\}$, $\{a_1, a_2\}$, $\ldots$, $\{a_1, \ldots, a_{49}\}$, $\ldots$
  - There must be some hidden relationships among frequent patterns!

# The Downward Closure Property of Frequent Patterns

□ Observation: From $TDB_1$: $T_1$: $\{a_1, \ldots, a_{50}\}$; $T_2$: $\{a_1, \ldots, a_{100}\}$

- We get a frequent itemset: $\{a_1, \ldots, a_{50}\}$
- Also, its subsets are all frequent: $\{a_1\}, \{a_2\}, \ldots, \{a_{50}\}, \{a_1, a_2\}, \ldots, \{a_1, \ldots, a_{49}\}, \ldots$
- There must be some hidden relationships among frequent patterns!

□ The downward closure (also called "Apriori") property of frequent patterns

- If **{beer, diaper, nuts}** is frequent, so is **{beer, diaper}**
- Every transaction containing {beer, diaper, nuts} also contains {beer, diaper}
- Apriori: Any subset of a frequent itemset must be frequent

A sharp knife for pruning!

# The Downward Closure Property of Frequent Patterns

- Observation: From $TDB_1$: $T_1$: $\{a_1, \ldots, a_{50}\}$; $T_2$: $\{a_1, \ldots, a_{100}\}$
  - We get a frequent itemset: $\{a_1, \ldots, a_{50}\}$
  - Also, its subsets are all frequent: $\{a_1\}$, $\{a_2\}$, …, $\{a_{50}\}$, $\{a_1, a_2\}$, …, $\{a_1, \ldots, a_{49}\}$, …
  - There must be some hidden relationships among frequent patterns!
- The <span style="color:red">downward closure (also called "Apriori")</span> property of frequent patterns
  - If **{beer, diaper, nuts}** is frequent, so is **{beer, diaper}**
  - Every transaction containing {beer, diaper, nuts} also contains {beer, diaper}
  - <span style="color:red">Apriori: Any subset of a frequent itemset must be frequent</span>
- Efficient mining methodology                          A sharp knife for pruning!
  - If <span style="color:red">any subset of an itemset S</span> is infrequent, then there is no chance for S to be frequent—why do we even have to consider S ?!

45

# Apriori Pruning and Scalable Mining Methods

- <u>Apriori pruning principle</u>: If there is any itemset which is infrequent, its superset should not even be generated!

  - (Agrawal & Srikant @VLDB'94, Mannila, et al. @ KDD' 94)

- Scalable mining Methods:  Three major approaches

  - Level-wise, join-based approach:

    - Apriori (Agrawal & Srikant@VLDB'94)

  - Vertical data format approach:

    - Eclat (Zaki, Parthasarathy, Ogihara, Li @KDD'97)

  - Frequent pattern projection and growth:

    - FPgrowth (Han, Pei, Yin @SIGMOD'00)

# Apriori: A Candidate Generation & Test Approach

- Outline of Apriori (level-wise, candidate generation and test)
    - Initially, scan DB once to get frequent 1-itemset
    - Repeat
        - Generate length-(k+1) candidate itemsets from length-k frequent itemsets
        - Test the candidates against DB to find frequent (k+1)-itemsets
        - Set k := k +1
    - Until no frequent or candidate set can be generated
    - Return all the frequent itemsets derived

# The Apriori Algorithm (Pseudo-Code)

$C_k$: Candidate itemset of size k

$F_k$ : Frequent itemset of size k

K := 1;
$F_k$ := {frequent items};   // frequent 1-itemset
**While** ($F_k$ != $\varnothing$) **do {**     // when $F_k$ is non-empty
    $C_{k+1}$ := candidates generated from $F_k$;  // candidate generation
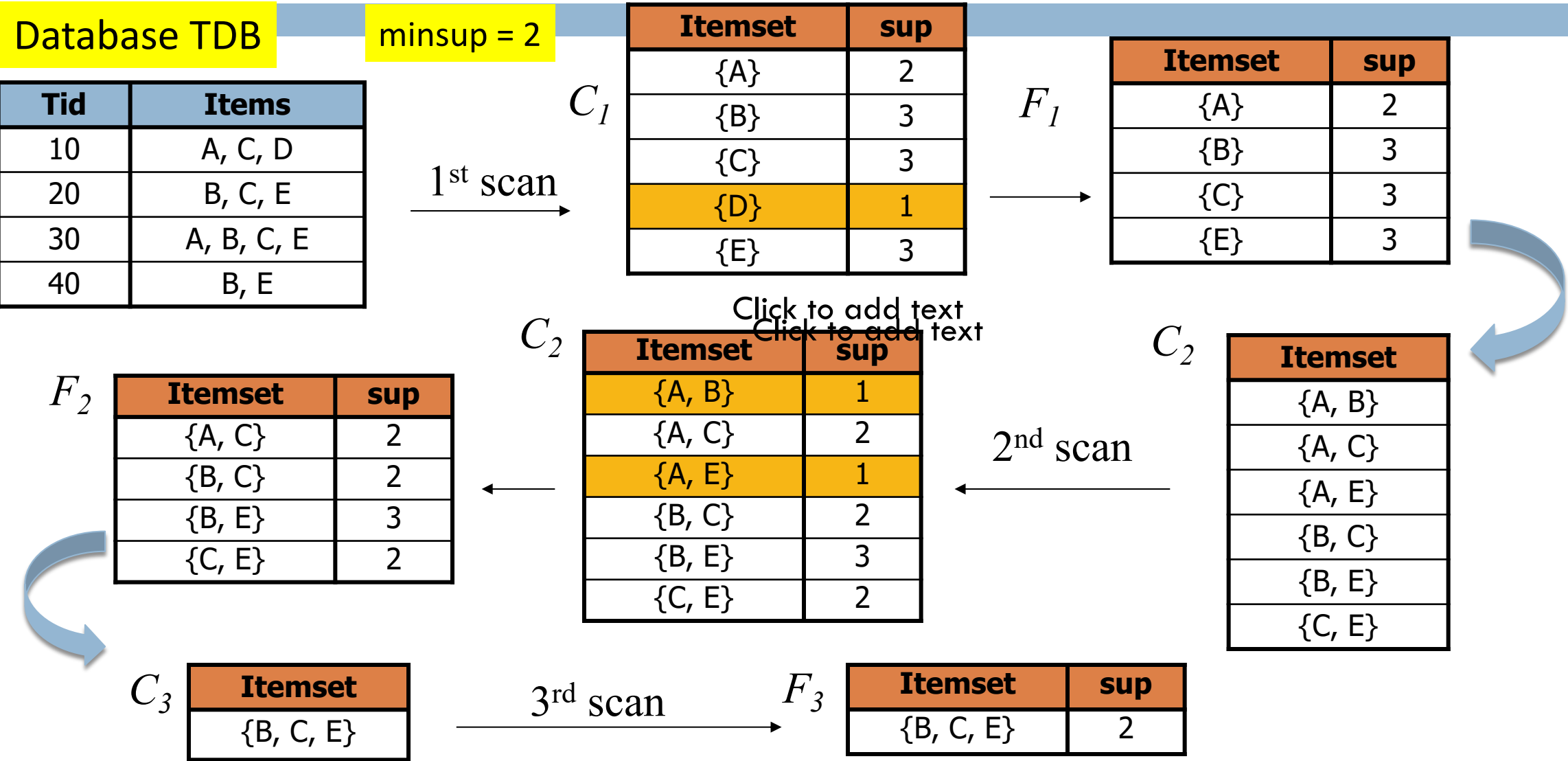    Derive $F_{k+1}$ by counting candidates in $C_{k+1}$ with respect to *TDB* at minsup;
    k := k + 1
    **}**
**return** $\cup_k F_k$             // return $F_k$ generated at each level

# The Apriori Algorithm—An Example

**Database TDB**   minsup = 2

| Tid | Items |
|---|---|
| 10 | A, C, D |
| 20 | B, C, E |
| 30 | A, B, C, E |
| 40 | B, E |

1st scan

$C_1$

| Itemset | sup |
|---|---|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {D} | 1 |
| {E} | 3 |

$F_1$

| Itemset | sup |
|---|---|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {E} | 3 |

Click to add text
Click to add text

$C_2$

| Itemset | sup |
|---|---|
| {A, B} | 1 |
| {A, C} | 2 |
| {A, E} | 1 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

2nd scan

$C_2$

| Itemset |
|---|
| {A, B} |
| {A, C} |
| {A, E} |
| {B, C} |
| {B, E} |
| {C, E} |

$F_2$

| Itemset | sup |
|---|---|
| {A, C} | 2 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$C_3$

| Itemset |
|---|
| {B, C, E} |

3rd scan

$F_3$

| Itemset | sup |
|---|---|
| {B, C, E} | 2 |

49

# The Apriori Algorithm—An Example

**Database TDB**  minsup = 2

| Tid | Items |
|-----|-------|
| 10 | A, C, D |
| 20 | B, C, E |
| 30 | A, B, C, E |
| 40 | B, E |

1st scan →

$C_1$

| Itemset | sup |
|---------|-----|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {D} | 1 |
| {E} | 3 |

$F_1$

| Itemset | sup |
|---------|-----|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {E} | 3 |

$C_2$

| Itemset |
|---------|
| {A, B} |
| {A, C} |
| {A, E} |
| {B, C} |
| {B, E} |
| {C, E} |

2nd scan

$C_2$

| Itemset | sup |
|---------|-----|
| {A, B} | 1 |
| {A, C} | 2 |
| {A, E} | 1 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$F_2$

| Itemset | sup |
|---------|-----|
| {A, C} | 2 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$C_3$

| Itemset |
|---------|
| {B, C, E} |

3rd scan

$F_3$

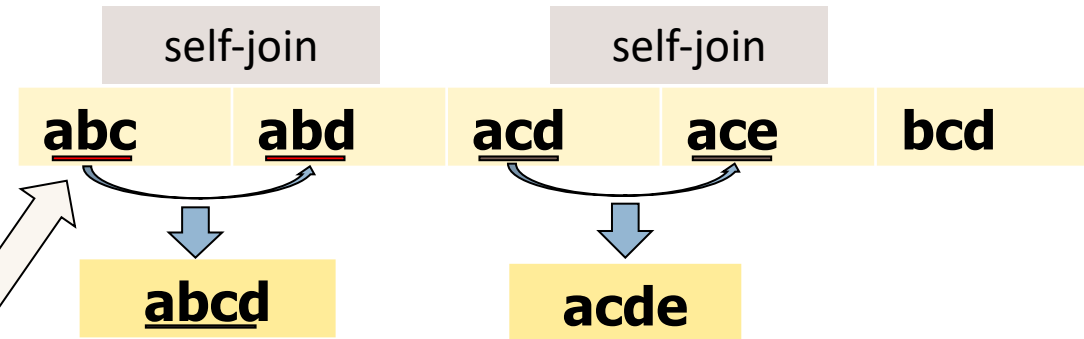| Itemset | sup |
|---------|-----|
| {B, C, E} | 2 |

Why?

50

# Apriori: Implementation Tricks

- How to generate candidates?
  - Step 1: self-joining $F_k$
  - Step 2: pruning

# Apriori: Implementation Tricks

- □ How to generate candidates?
  - ❑ Step 1: self-joining $F_k$
  - ❑ Step 2: pruning
- □ Example of candidate-generation
  - ❑ $F_3 = \{abc, abd, acd, ace, bcd\}$
  - ❑ Self-joining: $F_3*F_3$
    - ▪ *abcd* from *abc* and *abd*
    - ▪ *acde* from *acd* and *ace*

| self-join | | self-join | | |
|---|---|---|---|---|
| **abc** | **abd** | **acd** | **ace** | **bcd** |

**abcd**          **acde**

# Apriori: Implementation Tricks

- ▢ How to generate candidates?
  - ▢ Step 1: self-joining $F_k$
  - ▢ Step 2: pruning
- ▢ Example of candidate-generation
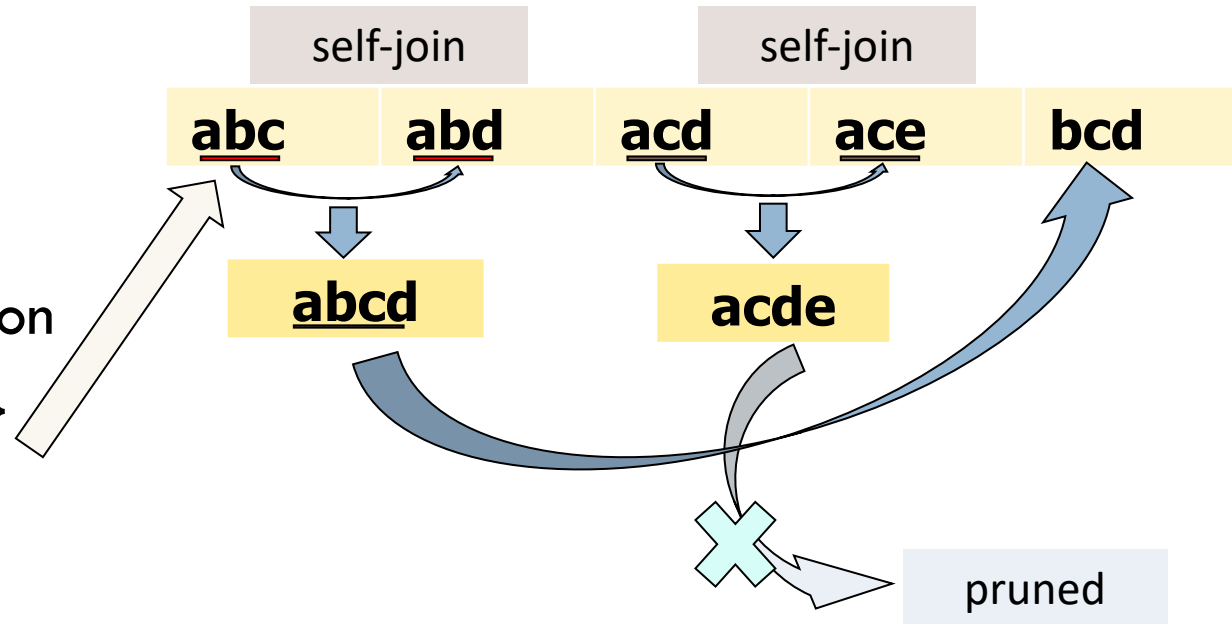  - ▢ $F_3 = \{abc, abd, acd, ace, bcd\}$
  - ▢ Self-joining: $F_3*F_3$
    - ■ *abcd* from *abc* and *abd*
    - ■ *acde* from *acd* and *ace*
  - ▢ Pruning:
    - ■ *acde* is removed because *ade* is not in $F_3$
  - ▢ $C_4 = \{abcd\}$

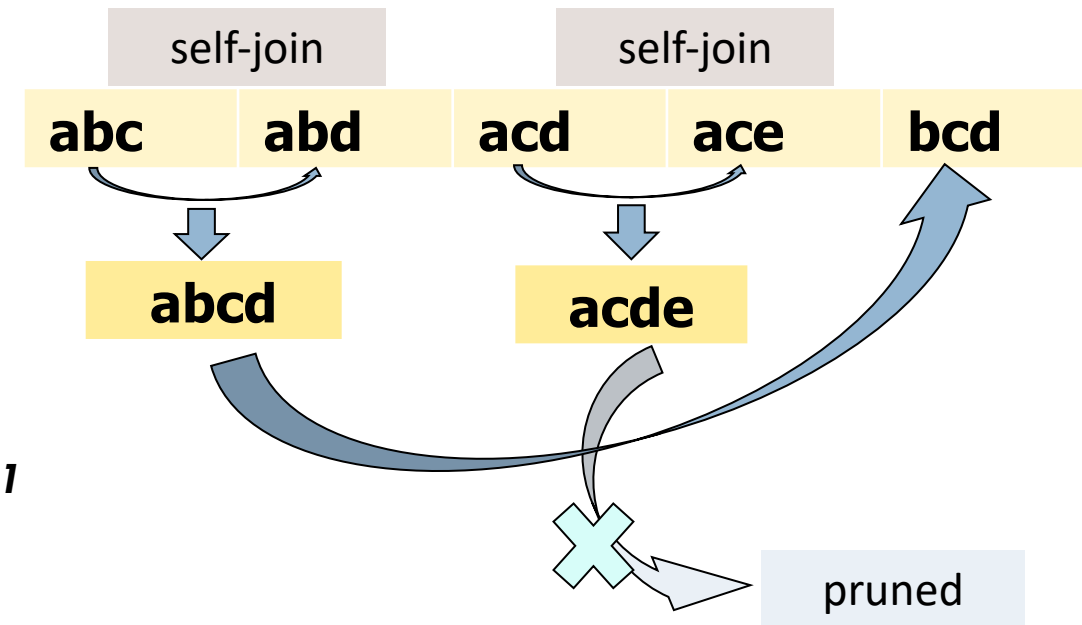| self-join | | self-join | | |
|---|---|---|---|---|
| **abc** | **abd** | **acd** | **ace** | **bcd** |

**abcd**        **acde**

pruned

# Candidate Generation: An SQL Implementation

☐ Suppose the items in $F_{k-1}$ are listed in an order

☐ Step 1: self-joining $F_{k-1}$

insert into $C_k$

select $p.item_1, p.item_2, …, p.item_{k-1}, q.item_{k-1}$

from $F_{k-1}$ as $p$, $F_{k-1}$ as $q$

where $p.item_1 = q.item_1, …, p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$

☐ Step 2: pruning

for all *itemsets c in $C_k$* do

for all *(k-1)-subsets s of c* do

if *(s is not in $F_{k-1}$)* then **delete** c **from** $C_k$

| self-join | | self-join | | |
|---|---|---|---|---|
| **abc** | **abd** | **acd** | **ace** | **bcd** |

**abcd**

**acde**

pruned

# Apriori Adv/Disadv

- *Advantages:*
  - Uses large itemset property
  - Easily parallelized
  - Easy to implement

- *Disadvantages:*
  - Assumes transaction database is memory resident
  - Requires up to m database scans

# Classification based on Association Rules (CBA)

- Why?
  - Can effectively uncover the correlation structure in data
  - AR are typically quite scalable in practice
  - Rules are often very intuitive
    - Hence classifier built on intuitive rules is easier to interpret

- When to use?
  - On large dynamic datasets where class labels are available and the correlation structure is unknown.
  - Multi-class categorization problems
  - E.g. Web/Text Categorization, Network Intrusion Detection