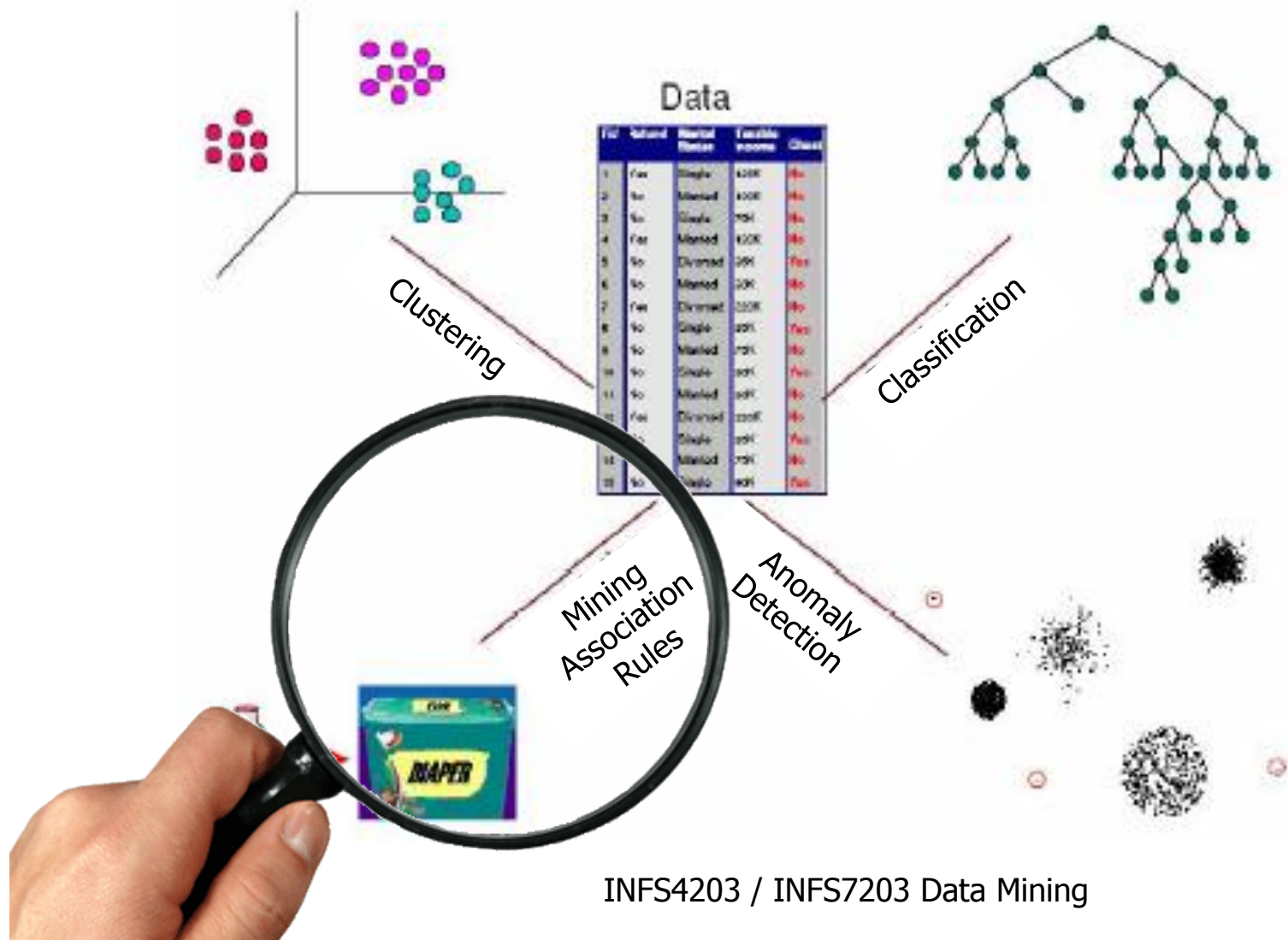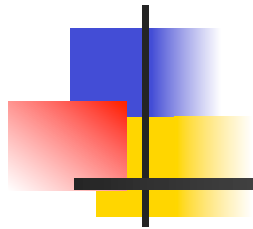# Data Mining Tasks

# Data Mining:
## Mining Association Rules

# What is Association Rule Mining?

- **Association rules mining**
  - Discovering interesting relations between objects in large databases

- **Market basket analysis**
  - The problem is to analyse customer buying habits by finding associations between the different items that customers place in their "shopping baskets"

    - 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

# Example: Market Basket Analysis

**Bread, Coke, Milk**

**Beer, Bread, Milk**

**Beer, Coke, Diapers, Milk**

**Beer, Bread, Diapers, Milk**

**Diapers, Milk**

Anything interesting?

Bread ⟹ Milk (100%)
Diapers ⟹ Beer (66%)
Diapers ⟹ Milk (100%)

Customers who buy diapers also tend to buy beer

⬇

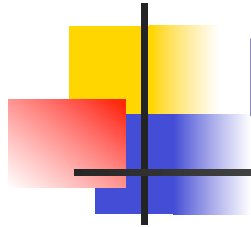Indentify potential cross-selling opportunities among related items

# Motivation (market basket analysis)

- If customers are buying milk, how <u>likely</u> is that they also buy bread?

- Such rules help retailers to:
  - <u>Plan the shelf space</u>: by placing milk close to bread they may increase the sales
  - <u>Provide advertisements/recommendation</u> to customers that are likely to buy some products
  - <u>Offer discounts</u> on items that are likely to be bought together to increase the sales

# Problem Statement

- **Given:**

  1. A database of transactions

  2. Each transaction is a list of items

     E.g.: items purchased by a customer in a visit

- **Find:**

  - All rules that correlate the presence of one set of items with another set of items

    E.g., 80% of customers who buy {diapers} tend to buy {beer, milk}.

# Association Rule Mining
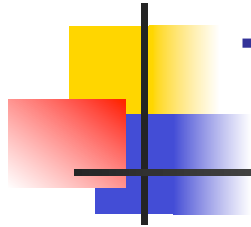
## Market-Basket transactions

| 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} → {Beer},
{Milk, Bread} → {Eggs,Coke},
{Beer, Bread} → {Milk},

...

...

# Transaction data: documents

- A text document dataset
  - Each document is treated as a **bag** of keywords

  doc1:        Student, Teach, School
  doc2:        Student, School
  doc3:        Teach, School, City, Game
  doc4:        Baseball, Basketball
  doc5:        Basketball, Player, Spectator
  doc6:        Baseball, Coach, Game, Team
  doc7:        Basketball, Team, City, Game

- Unusual words appearing together in a large number of documents, e.g., "Brad" and "Angelina," may indicate an interesting relationship.

# Formal Notations

- An **item**:  an item in a basket

- An **itemset** is a set of items.

  - E.g., X = {milk, bread, cereal} is an itemset.

- A **k-itemset** is an itemset with k items.

  - E.g., {milk, bread, cereal} is a 3-itemset

- A **transaction**: items purchased in a basket

  - it may have TID (transaction ID)

- A **transactional dataset**: A set of transactions

# Formal Notations

- An **association rule** is an implication of the form:
  $$X \rightarrow Y, \text{ where:}$$

  - $X, Y \subset I,$
  - $I$ is the set of all items,
  - $X \cap Y = \varnothing$

  e.g.: Bread, Butter → Milk

  Milk, Diapers → Beer

- Are all rules equally interesting?!

# Support and Confidence

- Rule X→Y

- **Rule support**

  - **(absolute) support** or support count:

    - frequency or **count** of an itemset $X \cup Y$
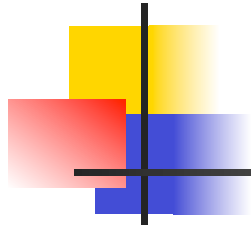
  - **(relative) support:**                                    P for probability.

    - probability that a transaction contains $X \cup Y$

    - support $(X→Y) = P(X \cup Y)$

  - **Support S:**

    - percentage of transactions that contain both X and Y

# Example: Frequent Itemsets

- Items={milk, coke, pepsi, beer, juice}.

- Support threshold = 3 baskets.

$B_1$ = {m, c, b}          $B_2$ = {m, p, j}

$B_3$ = {m, b}          $B_4$ = {c, j}

$B_5$ = {m, p, b}          $B_6$ = {m, c, b, j}

$B_7$ = {c, b, j}          $B_8$ = {b, c}

- Frequent itemsets: {m}, {c}, {b}, {j},

{m,b}, {b,c}, {c,j}.

# Support and Confidence

- Rule X→Y

- **Rule confidence:**

  - **Conditional probability**

  - Confidence (X→Y) = P(Y | X) = P(X ∪ Y) / P(X)

  - **Confidence C:**

    - Percentage of transactions that contain Y given they contain X

    - Measures how often items in Y appear in transactions that contain X

  - Example:

    - Support (printer→ink) **=** support (ink→printer)

    - But, confidence (printer→ink) **≠** confidence (ink→printer)
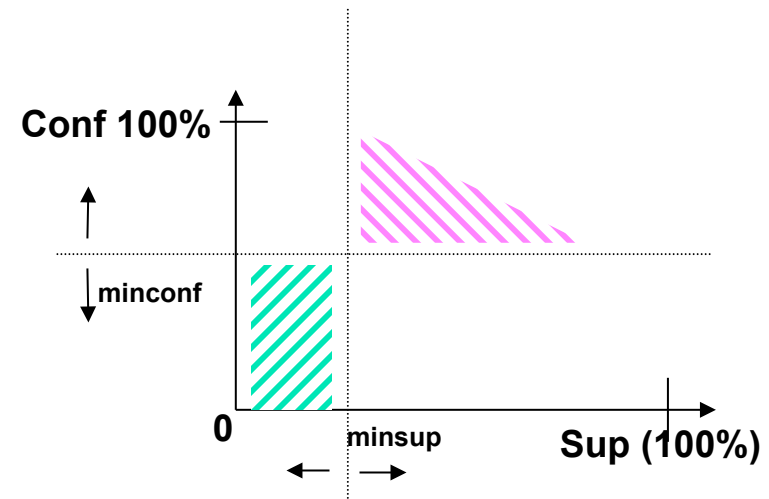
# Support and Confidence

- User-specified **thresholds:**
  - Minimum Support
    - *min_sup*
  - Minimum Confidence
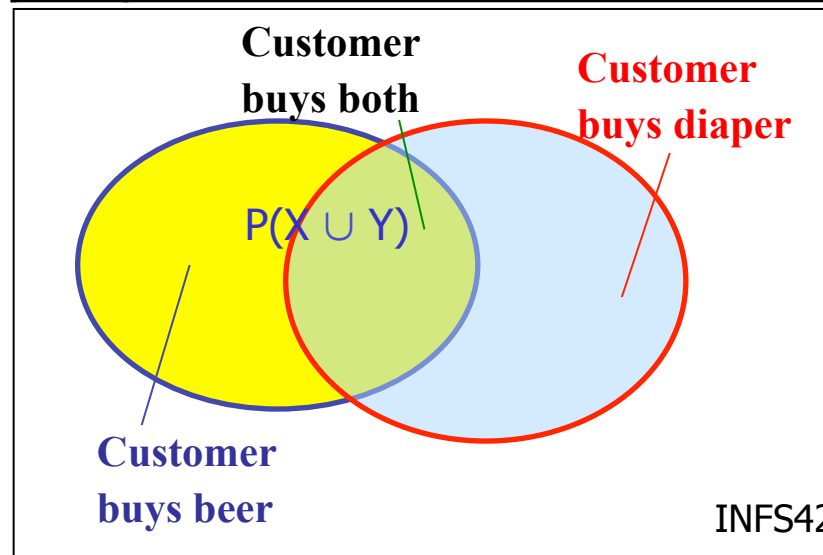    - *min_conf*



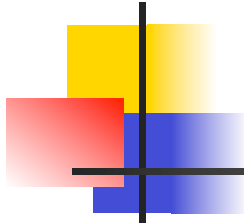Rules that satisfy both *min_sup* and *min_conf* are called **strong**

# Example

| 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 |

| itemset | support count | support |
|---------|--------------|---------|
| Milk | 2 | 40% |
| Beer | | |
| Nuts | | |
| Diaper | | |
| Eggs | | |
| Beer, Diaper | | |
| ... ... | | |

min_sup=50%



**Customer buys both**

**Customer buys diaper**

$P(X \cup Y)$

**Customer buys beer**

# Example

Find all the rules $X \rightarrow Y$ with minimum support and confidence

| 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 |

min_sup=50%
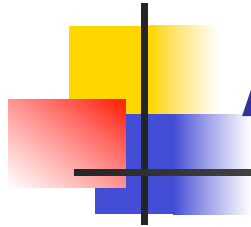min_conf = 50%

**Rule: Diaper→Beer**
support: 3 (60%)
confidence: P(Beer,Diaper)/**P(Diaper)**
= 60%/80% = **75%**

**Rule: Beer→Diaper**
support: 3 (60%)
confidence: P(Beer,Diaper)/**P(Beer)**
= 60%/60% = **100%**

# Association Rule Mining Task

- Given a set of transactions T, the goal of association rule mining is to find all rules having
    1. support ≥ *min_sup* threshold
    2. confidence ≥ *min_conf* threshold

- Brute-force approach:
    - List all possible association rules
    - Compute the support and confidence for each rule
    - Prune rules that fail the *min_sup* or *min_conf* thresholds
    - ⇒ Computationally prohibitive!

# Mining Association Rules

1. **Apriori Algorithm**

2. Frequent Pattern (FP) Growth Algorithm
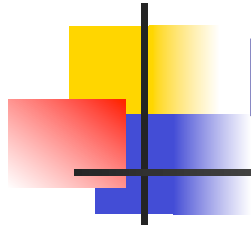
# Mining Association Rules

| 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} → {Beer} (s=0.4, c=0.67)
{Milk,Beer} → {Diaper} (s=0.4, c=1.0)
{Diaper,Beer} → {Milk} (s=0.4, c=0.67)
{Beer} → {Milk,Diaper} (s=0.4, c=0.67)
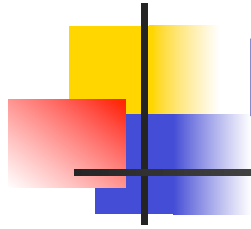{Diaper} → {Milk,Beer} (s=0.4, c=0.5)
{Milk} → {Diaper,Beer} (s=0.4, c=0.5)

## Observations:

• All the above rules are partitions of the same itemset:
   {Milk, Diaper, Beer}

• Rules originating from the same itemset have identical support but can have different confidence

• Thus, **decouple** the support and confidence requirements

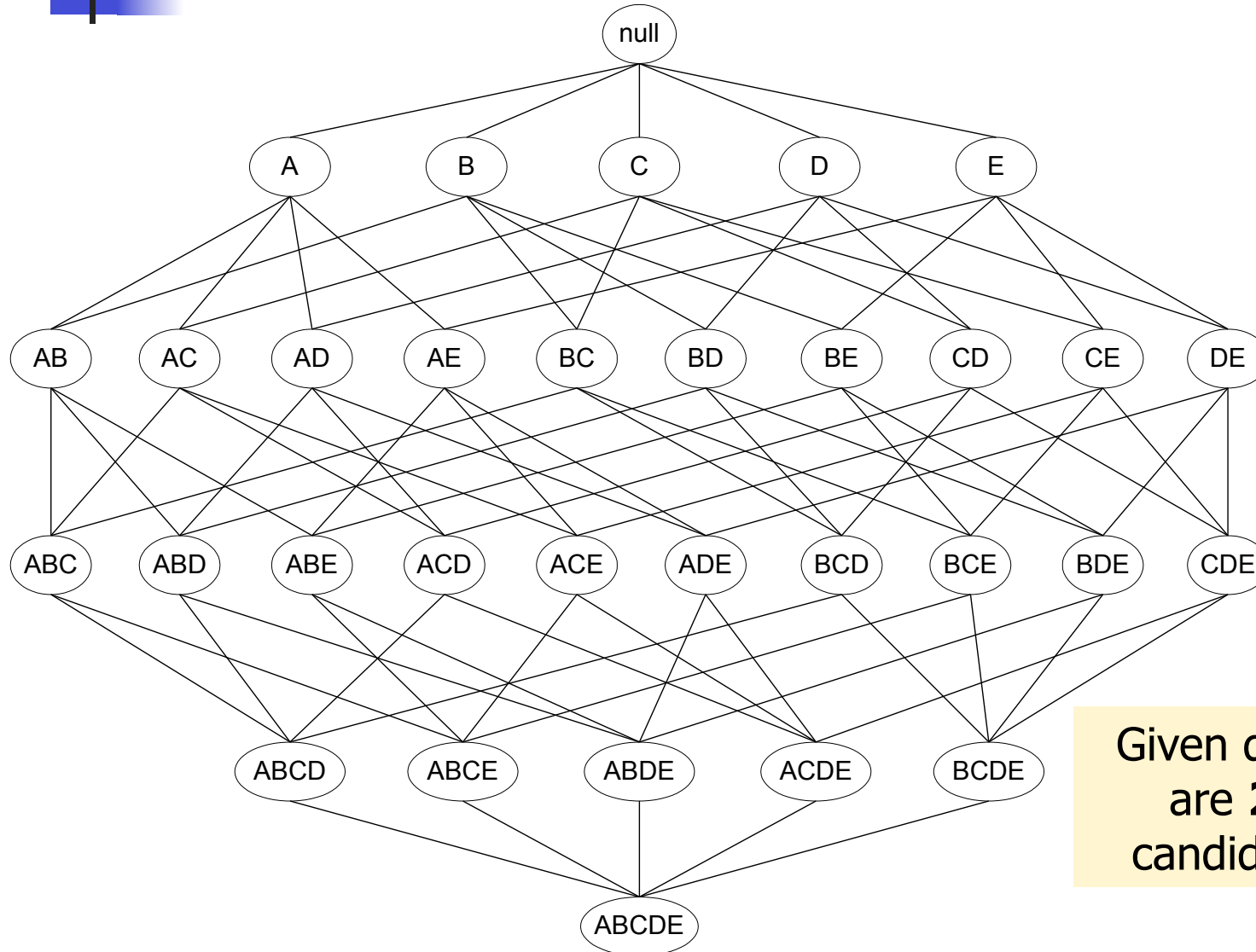# Mining Association Rules

- Two-step approach:

    1. **Frequent Itemset Generation**
        - Generate all itemsets whose support ≥ minsup

    2. **Rule Generation**
        - Generate high confidence rules from each frequent itemset
            - each rule is a partitioning of a frequent itemset

# 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 candidate frequent itemset
  - Count the support of each candidate by scanning the database

**Transactions**

**List of Candidates**

| 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

M

w

  - Match each transaction <u>against every candidate</u>
  - Complexity ~ O(NMw) => Expensive since M = $2^d$ !!!

# Scale of the Problem

- ## WalMart:
  - Sells 100,000 items and can store billions of baskets.

- ## The Web
  - Has  billions of words and many billions of pages.

# 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



null

A  B  C  D  E

AB  AC  AD  AE  BC  BD  BE  CD  CE  DE

ABC  ABD  ABE  ACD  ACE  ADE  BCD  BCE  BDE  CDE

ABCD  ABCE  ABDE  ACDE  BCDE

ABCDE

Found to be
Infrequent

if AB is not frequent,
ABC must be not frequent, but we cannot
assume A is not frequent.

Pruned
supersets

# The Apriori Algorithm

Begin

**Initialize** the *candidate* Itemsets with single items in database (K=1)

**Scan** the database and **count** the frequency of the candidate itemsets then *frequent* Itemsets are <u>decided</u> based on the user specified **min_sup**

Any **new** frequent Itemsets?

**NO**

**YES**

**Expand** the frequent itemsets with one more item to generate new candidate itemsets.

Stop

# The Apriori Algorithm—An Example

$Sup_{min} = 2$

Database TDB

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

1st scan

candidates

$C_1$

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

frequent items

$L_1$

| 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} |

$C_2$

2nd scan

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

$L_2$

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

$C_3$

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

3rd scan

$L_3$

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

# The Apriori Algorithm (Pseudo-Code)

$C_k$: Candidate itemset of size k
$L_k$ : frequent itemset of size k

$L_1$ = {frequent items};
**for** ($k$ = 1; $L_k$ !=$\varnothing$; $k$++) **do begin**
    $C_{k+1}$ = **candidates generated** from $L_k$;
    **for each** transaction $t$ in database do
        increment the count of all candidates in $C_{k+1}$ that
         are contained in $t$
    $L_{k+1}$  = candidates in $C_{k+1}$ with min_support
    **end**
**return** $\cup_k$ $L_k$;

# Implementation of Apriori

- How to generate candidates?
  - Step 1: self-joining $L_k$
  - Step 2: pruning
- Example of Candidate-generation
  - $L_3=\{abc,\ abd,\ acd,\ ace,\ bcd\}$
  - Self-joining: $L_3*L_3$

    开头两个一样的 组合在一起
    - *abcd* from *abc* and *abd*     abc and abd
    - *acde* from *acd* and *ace*     acd and ace
  - Pruning:     all aubset should be frequent. but in acde, there is not ade
    - *acde* is removed because *ade* is not in $L_3$
  - $C_4 = \{abcd\}$

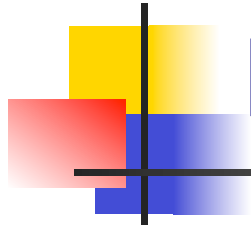# Candidate Generation: SQL Implementation

- SQL Implementation of candidate generation
  - Suppose the items in $L_{k-1}$ are listed in an order

  - Step 1: self-joining $L_{k-1}$

    insert into $C_k$

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

    from $L_{k-1}\ p,\ L_{k-1}\ q$

    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

    forall *itemsets c in $C_k$* do

        forall *(k-1)-subsets s of c* do

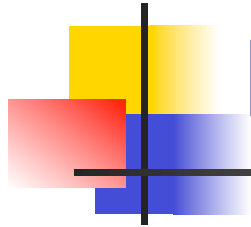            **if** *(s is not in $L_{k-1}$)* **then delete** *c* **from** *C*

# Rule Generation

- Given a frequent itemset L, find all non-empty subsets f ⊂ L such that f → L − f **satisfies** the minimum confidence requirement

  - If {A,B,C,D} is a frequent itemset, candidate rules:

    ABC →D,    ABD →C,    ACD →B, BCD →A,
    A →BCD,    B →ACD,    C →ABD,    D →ABC
    AB →CD,    AC → BD,    AD → BC,    BC→AD,
    BD →AC,    CD →AB,

# Factors Affecting Complexity

- ## Choice of minimum support threshold

  - lower support threshold results in more frequent itemsets

- ## Dimensionality (number of items) in the data set

  - 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

  - This may increase max length of frequent itemsets

# Midterm Guidelines

- Wednesday Sept 19th @8:00am
- Location:
  - 01−E109
  - 76-228

- Duration: 90 minutes

- Up to and including material covered last week!
- Problem solving, short fill-in questions
- Calculator & Student ID
- Please check last years final exam questions!