# Data Mining

# Knowledge Discovery in Databases
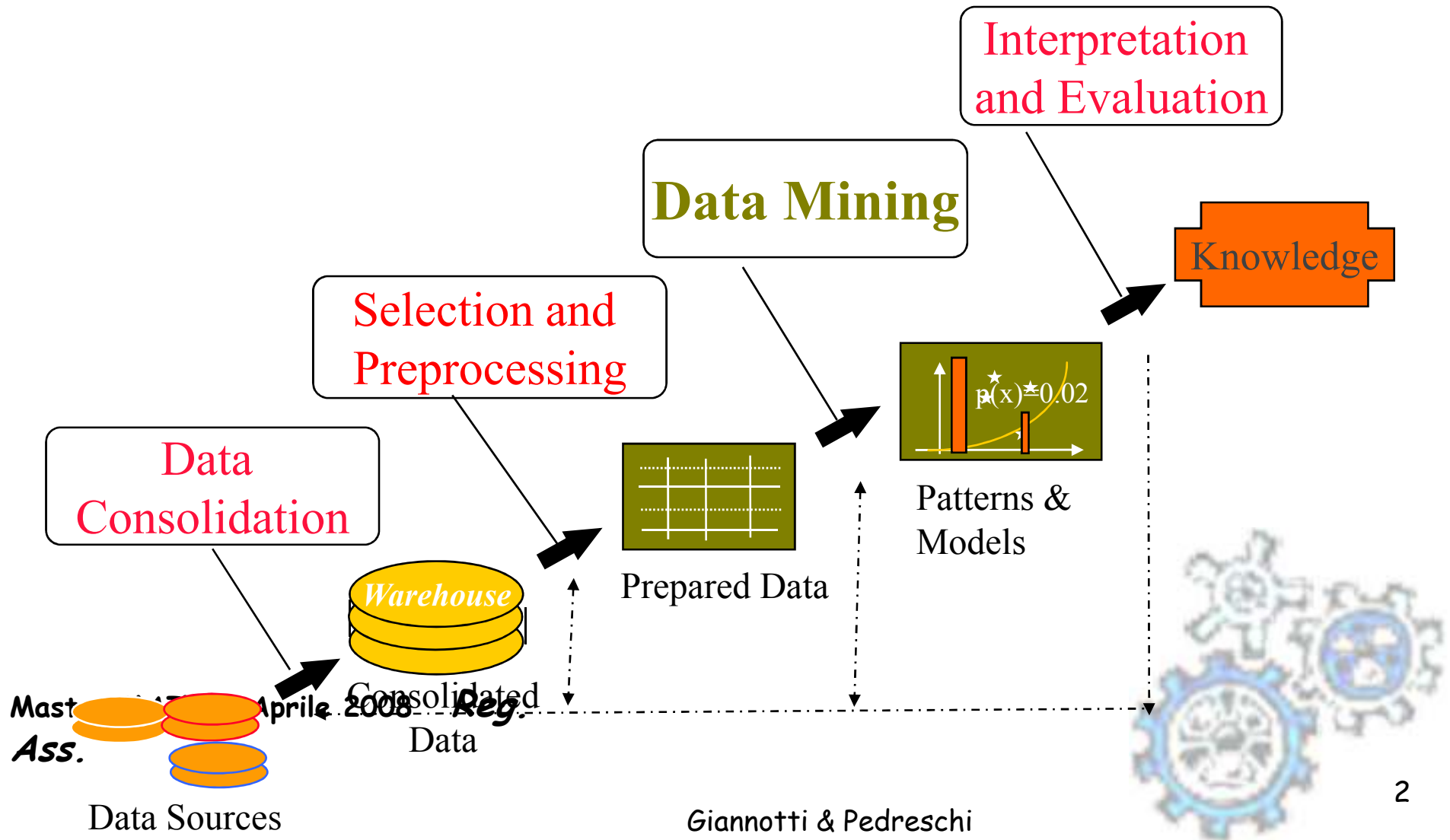
**Fosca Giannotti and Dino Pedreschi**

**Pisa KDD Lab, ISTI-CNR & Univ. Pisa**

**http://www-kdd.cnuce.cnr.it/**

**MAINS – Master in Management dell'Innovazione**
**Scuola Superiore S. Anna**

# KDD Process

Interpretation and Evaluation

Data Mining

Knowledge

Selection and Preprocessing

$p(x) = 0.02$

Data Consolidation

Patterns & Models

Warehouse

Prepared Data

Master ... Aprile 2008 ... Reg.

Ass.

Consolidated Data

Data Sources

# Association rules and market basket analysis

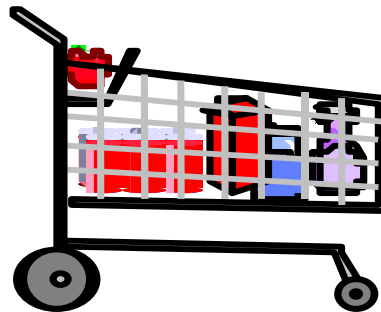# Market Basket Analysis: the context

Customer buying habits by finding associations and correlations between the different items that customers place in their "shopping basket"

Milk, eggs, sugar, bread

Milk, eggs, cereal, bread

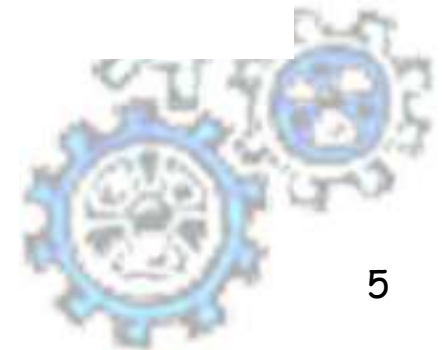Eggs, sugar

Customer1

Customer2

Customer3

Giannotti & Pedreschi
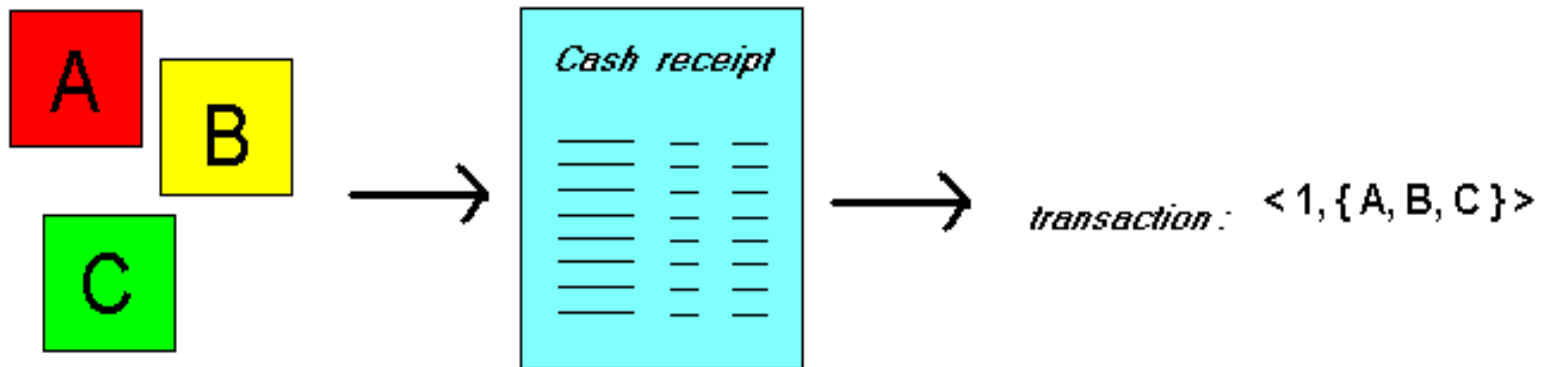
# Market Basket Analysis: the context

Given: a database of customer **transactions**, where each transaction is a **set of items**

☒ Find groups of items which are **frequently purchased together**

Giannotti & Pedreschi

# Goal of MBA

- **Extract information on purchasing behavior**
- **Actionable information: can suggest**
  - new store layouts
  - new product assortments
  - which products to put on promotion
- **MBA applicable whenever a customer purchases multiple things in proximity**
  - credit cards
  - services of telecommunication companies
  - banking services
  - medical treatments

# MBA: applicable to many other contexts

**Telecommunication:**

Each customer is a transaction containing the set of customer's phone calls
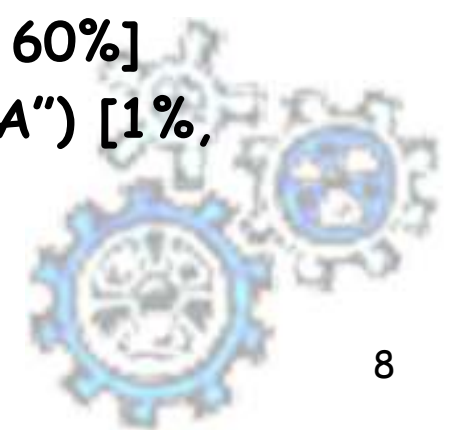
**Atmospheric phenomena:**

Each time interval (e.g. a day) is a transaction containing the set of observed event (rains, wind, etc.)

**Etc.**

# Association Rules

- Express how product/services relate to each other, and tend to group together

- "if a customer purchases three-way calling, then will also purchase call-waiting"

- simple to understand

- actionable information: bundle three-way calling and call-waiting in a single package

- Examples.
  - Rule form: "Body $\rightarrow$ Head [support, confidence]".
  - buys(x, "diapers") $\rightarrow$ buys(x, "beers") [0.5%, 60%]
  - major(x, "CS") ^ takes(x, "DB") $\rightarrow$ grade(x, "A") [1%, 75%]

Giannotti & Pedreschi

# Useful, trivial, unexplicable

- **Useful**: "On Thursdays, grocery store consumers often purchase diapers and beer together".

- **Trivial**: "Customers who purchase maintenance agreements are very likely to purchase large appliances".

- **Unexplicable**: "When a new hardaware store opens, one of the most sold items is toilet rings."

Giannotti & Pedreschi

# 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!

Giannotti & Pedreschi

# 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} $\rightarrow$ {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:

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

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

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

Giannotti & Pedreschi

12

# 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**!
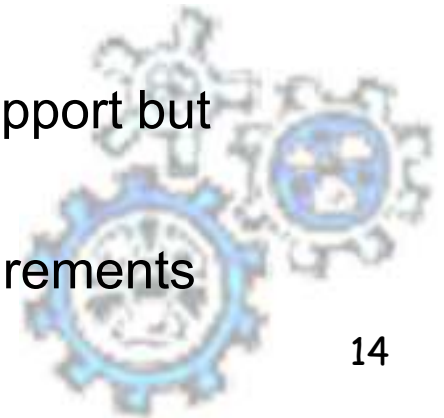
# Mining Association Rules

| 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

Giannotti & Pedreschi

# 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**

# Basic Apriori Algorithm

## Problem Decomposition

⏱ **Find the *frequent itemsets*: the sets of items that satisfy the support constraint**

- ◆ **A subset of a frequent itemset is also a frequent itemset**, i.e., if {*A,B*} is a frequent itemset, both {*A*} and {*B*} should be a frequent itemset

- ◆ Iteratively find frequent itemsets with cardinality from 1 to *k (k*-itemset*)*

⏱ **Use the frequent itemsets to generate association rules.**

Giannotti & Pedreschi

# Frequent Itemset Generation



**Given d items, there are 2<sup>d</sup> possible candidate itemsets**

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

Master MAINS, Aprile 2008    *Reg.*
*Ass.*

Giannotti & Pedreschi

17

# Reducing Number of Candidates

- **Apriori principle**:
  - If an itemset is frequent, then all of its subsets must also be frequent

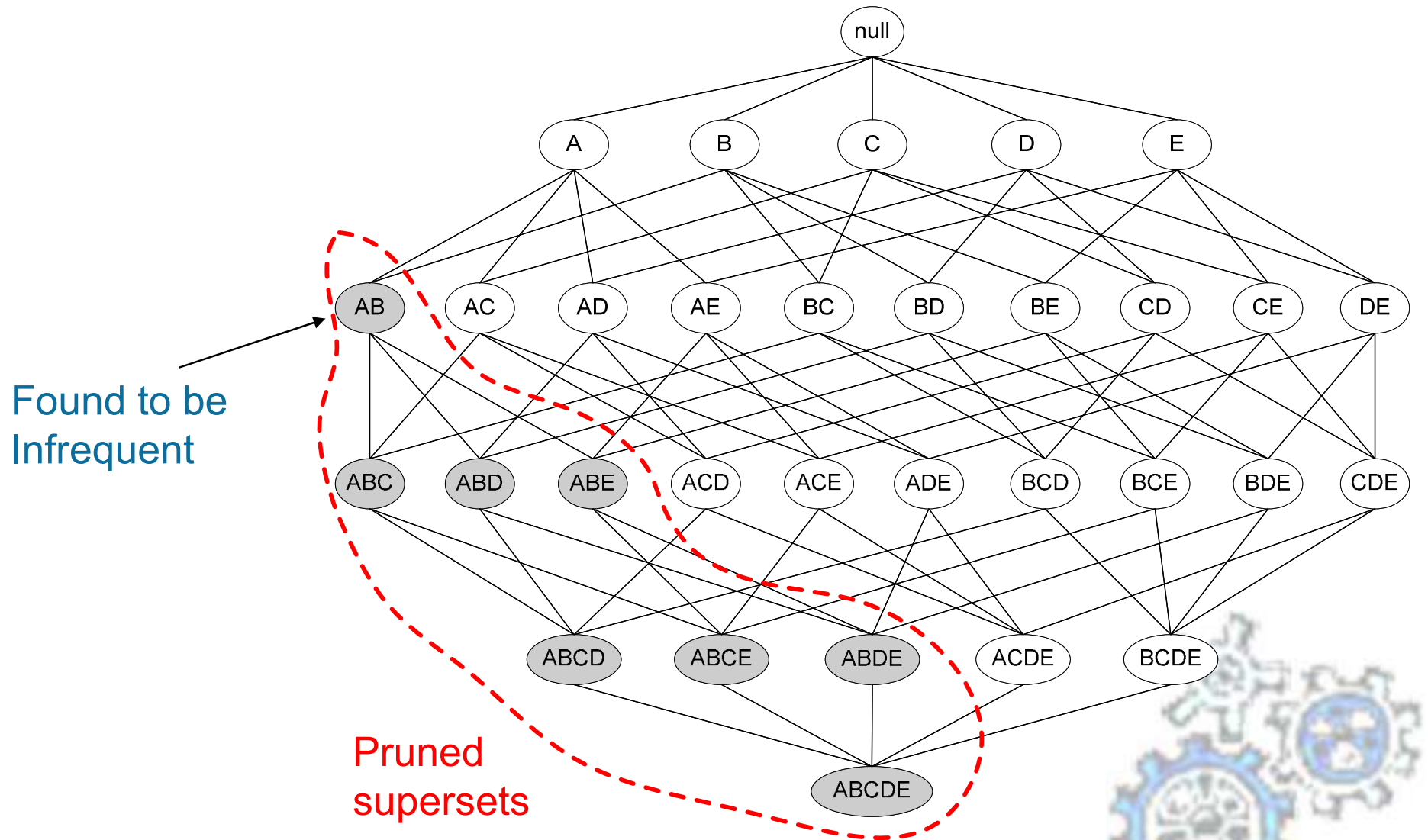- **Apriori principle holds due to the following property of the support measure:**

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

  - 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

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

Items (1-itemsets)

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

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Triplets (3-itemsets)

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

Minimum Support = 3

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

Giannotti & Pedreschi

# *Apriori Execution Example* *(min_sup = 2)*

Database TDB

| TID | Items |
|-----|-------|
| 100 | 1 3 4 |
| 200 | 2 3 5 |
| 300 | 1 2 3 5 |
| 400 | 2 5 |

Scan TDB →

$C_1$

| itemset | sup. |
|---------|------|
| {1} | 2 |
| {2} | 3 |
| {3} | 3 |
| {4} | 1 |
| {5} | 3 |

→

$L_1$

| itemset | sup. |
|---------|------|
| {1} | 2 |
| {2} | 3 |
| {3} | 3 |
| {5} | 3 |

$L_2$

| itemset | sup |
|---------|-----|
| {1 3} | 2 |
| {2 3} | 2 |
| {2 5} | 3 |
| {3 5} | 2 |

← $C_2$

| itemset | sup |
|---------|-----|
| {1 2} | 1 |
| {1 3} | 2 |
| {1 5} | 1 |
| {2 3} | 2 |
| {2 5} | 3 |
| {3 5} | 2 |

← Scan TDB ← $C_2$

| itemset |
|---------|
| {1 2} |
| {1 3} |
| {1 5} |
| {2 3} |
| {2 5} |
| {3 5} |

$C_3$

| itemset |
|---------|
| {2 3 5} |

Scan TDB →

$L_3$

| itemset | sup |
|---------|-----|
| {2 3 5} | 2 |

Giannotti & Pedreschi

# Multidimensional AR

Associations between values of different attributes :

| CID | nationality | age | income |
|-----|-------------|-----|--------|
| 1 | Italian | 50 | low |
| 2 | French | 40 | high |
| 3 | French | 30 | high |
| 4 | Italian | 50 | medium |
| 5 | Italian | 45 | high |
| 6 | French | 35 | high |

RULES:

**nationality** = French $\Rightarrow$ **income** = high [50%, 100%]

**income** = high $\Rightarrow$ **nationality** = French [50%, 75%]

**age** = 50 $\Rightarrow$ **nationality** = Italian [33%, 100%]

Giannotti & Pedreschi

# Hierarchy of concepts

Department — Sector — Family — Product

FoodStuff
- Frozen
- Refrigerated
- Fresh
  - Vegetable
  - Fruit
    - Banana
    - Apple
    - Orange
    - Etc...
  - Dairy
  - Etc....
- Bakery
- Etc...