

# Association Analysis

## IE500 Data Mining



# Outline

1. What is Association Analysis?
2. Frequent Itemset Generation
3. Rule Generation
4. Handling Continuous and Categorical Attributes
5. Interestingness Measures
6. Sequential Pattern Mining

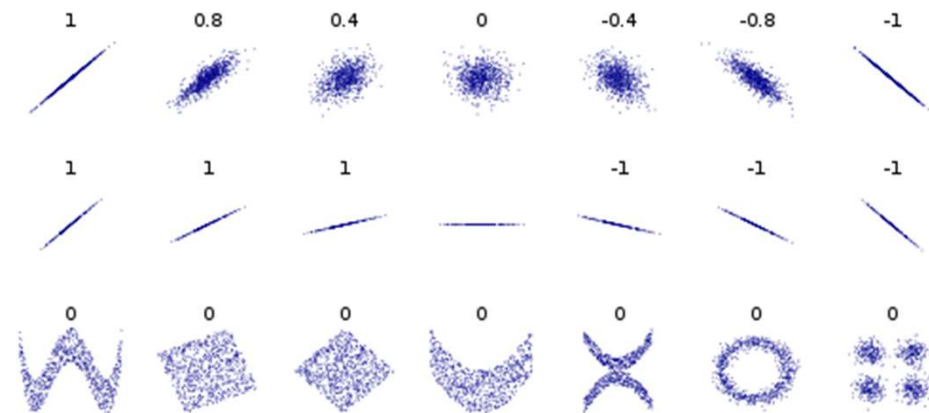
# Correlation Analysis

- Correlation analysis measures the degree of dependency between **two variables**
  - Continuous variables: Pearson's correlation coefficient (PCC)
  - Binary variables: Phi coefficient

$$PCC(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

$$Phi(x, y) = \frac{f_{11}f_{00} - f_{01}f_{10}}{\sqrt{f_{1+}f_{+1}f_{0+}f_{+0}}}$$

- Value range [-1,1]
  - 1 : positive correlation
  - 0 : variables independent
  - 1 : negative correlation



# Example: Correlations between Products in Shopping Baskets

	P1	P2	P3	P4	P5
Basket 1	1	1	0	1	1
Basket 2	1	0	0	1	1
Basket 3	1	0	0	0	1

1 : always bought together  
 0 : sometimes bought together  
 -1 : never bought together

Correlation Matrix (Correlation Matrix)

Table View  Pairwise Table  Plot View  Annotations

Attributes	ThinkPad X...	Asus EeePC	HP Laserjet...	2 GB DDR3...	8 GB DDR3...	Lenovo Tab...	Netbook-Sc...	HP CE50 T...	LT Laser M...	LT Minimaus
ThinkPad X2	1	-1	0.356	-0.816	0.612	0.583	-0.667	0.356	0.167	-0.408
Asus EeePC	-1	1	-0.356	0.816	-0.612	-0.583	0.667	-0.356	-0.167	0.408
HP Laserjet	0.356	-0.356	1	-0.218	-0.327	0.356	-0.535	1	-0.089	-0.655
2 GB DDR3	-0.816	0.816	-0.218	1	-0.500	-0.816	0.816	-0.218	0	0.200
8 GB DDR3	0.612	-0.612	-0.327	-0.500	1	0.102	-0.408	-0.327	0.102	0
Lenovo Tabl	0.583	-0.583	0.356	-0.816	0.102	1	-0.667	0.356	-0.250	0
Netbook-Scf	-0.667	0.667	-0.535	0.816	-0.408	-0.667	1	-0.535	0.167	0.408
HP CE50 To	0.356	-0.356	1	-0.218	-0.327	0.356	-0.535	1	-0.089	-0.655
LT Laser Ma	0.167	-0.167	-0.089	0	0.102	-0.250	0.167	-0.089	1	-0.408
LT Minimaus	-0.408	0.408	-0.655	0.200	0	0	0.408	-0.655	-0.408	1

**Shortcoming:** Measures correlation only between two items but not between multiple items, e.g. {ThinkPad, Cover} → {Minimaus}

# 1. Association Analysis

- Association analysis can find **multiple item co-occurrence relationships** (descriptive method)
- first algorithms developed in the early 90s at IBM by Agrawal & Srikant
- initially used for **shopping basket analysis** to find out how items purchased are related
- later extended to more complex data structures
  - sequential patterns
  - subgraph patterns
- and other application domains
  - web usage mining, social science, life science



# Association Analysis

- 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.
- Examples of **Association Rules**
  - {Diaper} → {Beer}
  - {Beer, Bread} → {Milk}
  - {Milk, Bread} → {Eggs, Coke}

Implication denotes  
co-occurrence, not causality!

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

Shopping Transactions

# Definition: Frequent Itemset

- **Itemset**
  - 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.  $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$
- **Support (s)**
  - Fraction of transactions that contain an itemset
  - e.g.  $s(\{\text{Milk, Bread, Diaper}\}) = 2/5 = 0.4$
- **Frequent Itemset**
  - An itemset whose support is greater than or equal to a minimal support (minsup) threshold specified by the user

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

Shopping Transactions

# Definition: Association Rule

- Association Rule
  - An implication of the form  $X \rightarrow Y$ , where  $X$  and  $Y$  are itemsets
  - Interpretation: when  $X$  occurs,  $Y$  occurs with a certain probability

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

Shopping Transactions

- More formally, it's a *conditional probability*
  - $P(Y|X)$  given  $X$ , what is the probability of  $Y$ ?

# Definition: Association Rule

- Association Rule
  - Example:
 
$$\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$$

**Condition**
**Consequent**

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

Shopping Transactions

- Rule Evaluation Metrics
  - **Support s:**  
Fraction of total transactions which contain both X and Y

$$s(X \rightarrow Y) = \frac{|X \cup Y|}{|T|}$$

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

- **Confidence c:**  
Measures how often items in Y appear in transactions that contain X

$$c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$$

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

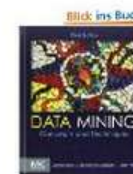
# Main Challenges concerning Association Analysis

1. Mining association rules from large amounts of data can be **computationally expensive**
  - algorithms need to apply smart pruning strategies
2. Algorithms often discover a **large number of associations**
  - many of them are uninteresting or redundant
  - the user needs to select the subset of the associations that is relevant given her task at hand

Wird oft zusammen gekauft



+



Preis für beide: EUR 138,00

Beides in den Einkaufswagen

Verfügbarkeit und Versanddetails anzeigen

- Dieser Artikel:** Introduction to Data Mining von Pang-Ning Tan Taschenbuch EUR 85,05
- Data Mining: Concepts and Techniques (Morgan Kaufmann Series in Data Management Systems)

# The Association Rule Mining Task

- Given a set of **transactions T**, the goal of association rule mining is to **find all rules** having
  - support  $\geq$  *minsup* threshold
  - confidence  $\geq$  *minconf* threshold
- ***minsup*** and ***minconf*** are provided by the user
- Brute Force Approach:
  1. List all possible association rules
  2. Compute the support and confidence for each rule
  3. Remove rules that fail the *minsup* and *minconf* thresholds

**⇒ Computationally prohibitive due to large number of candidates!**

# Mining Association Rules

- Example rules:  
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\} (s=0.4, c=0.67)$   
 $\{\text{Milk, Beer}\} \rightarrow \{\text{Diaper}\} (s=0.4, c=1.0)$   
 $\{\text{Diaper, Beer}\} \rightarrow \{\text{Milk}\} (s=0.4, c=0.67)$   
 $\{\text{Beer}\} \rightarrow \{\text{Milk, Diaper}\} (s=0.4, c=0.67)$   
 $\{\text{Diaper}\} \rightarrow \{\text{Milk, Beer}\} (s=0.4, c=0.5)$   
 $\{\text{Milk}\} \rightarrow \{\text{Diaper, Beer}\} (s=0.4, c=0.5)$

TID	Items
1	Bread, <b>Milk</b>
2	Bread, <b>Diaper, Beer</b> , Eggs
3	<b>Milk, Diaper, Beer</b> , Coke
4	Break, <b>Milk, Diaper, Beer</b>
5	Bread, <b>Milk, Diaper</b> , Coke

Shopping Transactions

- Observations:
  - All the above rules are binary partitions of the same itemset:  
 $\{\text{Milk, Diaper, Beer}\}$
  - Rules originating from the same itemset have identical support  $s$ 
    - but can have different confidence
  - Thus, we may decouple the support and confidence requirements

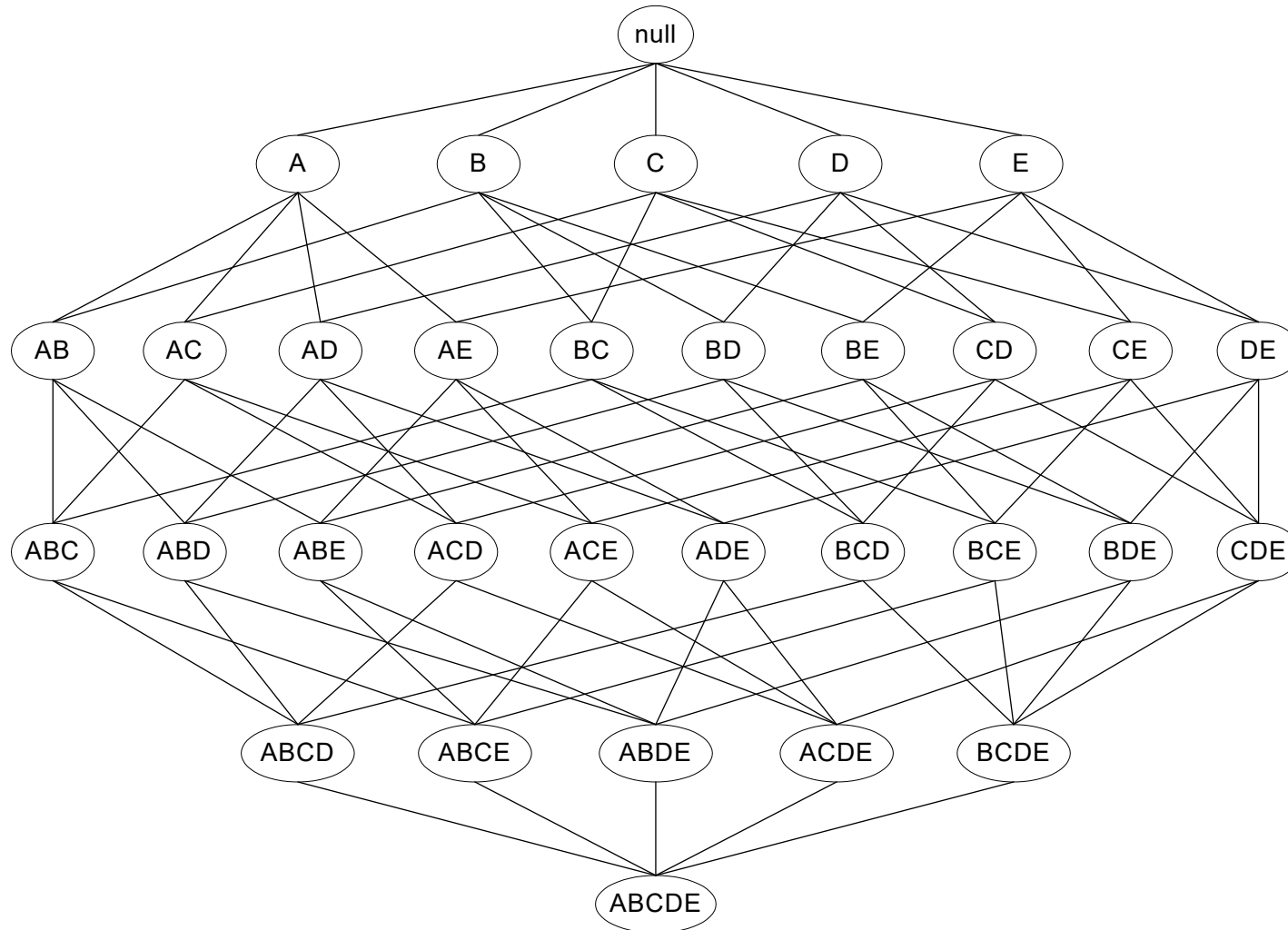
$$s(X \rightarrow Y) = \frac{|X \cup Y|}{|T|}$$

# Apriori Algorithm: Basic Idea

- 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

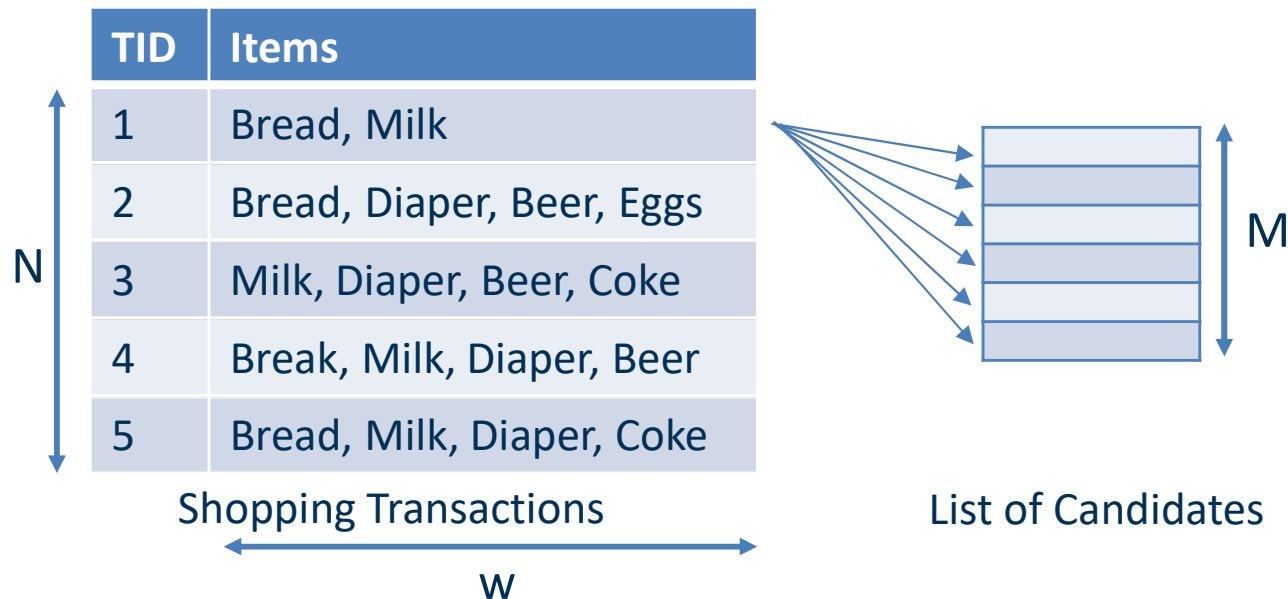
# 2. Frequent Itemset Generation

- Given d items, there are  $2^d$  candidate itemsets!



# Brute-force Approach

- Each itemset in the lattice is a **candidate** frequent itemset
- Count the support of each candidate by scanning the database  
= Match each transaction against every candidate



- Complexity  $\sim O(NMw) \rightarrow$  **Expensive since  $M = 2^d$**

# Brute-force Approach

- Amazon sells 12M different products (as of 2023)
  - That is  $2^{12.000.000}$  possible itemsets
    - That's a 3.6M digit number
  - Today's supercomputers: 1,200 Petaflops, i.e.,  $1.2 \times 10^{18}$  floating point operations per second
  - Even if an itemset could be checked with one single floating point operation this would take  $\sim 10^{3,612,334}$  years (age of universe:  $1.4 \times 10^{10}$  years)
- However:
  - Most itemsets will not be frequent at all, e.g., books on Chinese calligraphy, Inuit cooking, and data mining bought together
  - Thus, smarter algorithms should be possible
  - Intuition for smarter algorithms: All itemsets containing Inuit cooking are likely infrequent



# Anti-Monotonicity of Support

- What happens when an itemset gets larger?
  - $s(\{\text{Milk}\}) = 0.8$
  - $s(\{\text{Milk}, \text{Diaper}\}) = 0.6$
  - $s(\{\text{Milk}, \text{Diaper}, \text{Beer}\}) = 0.4$
  
  - $s(\{\text{Bread}\}) = 0.8$
  - $s(\{\text{Bread}, \text{Milk}\}) = 0.6$
  - $s(\{\text{Bread}, \text{Milk}, \text{Diaper}\}) = 0.4$
- There is a pattern here!

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Break, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

# Reducing the Number of Candidates

- There is a pattern here!
  - It is called anti-monotonicity of support
- If  $X$  is a subset of  $Y$ 
  - $s(Y)$  is at most as large as  $s(X)$

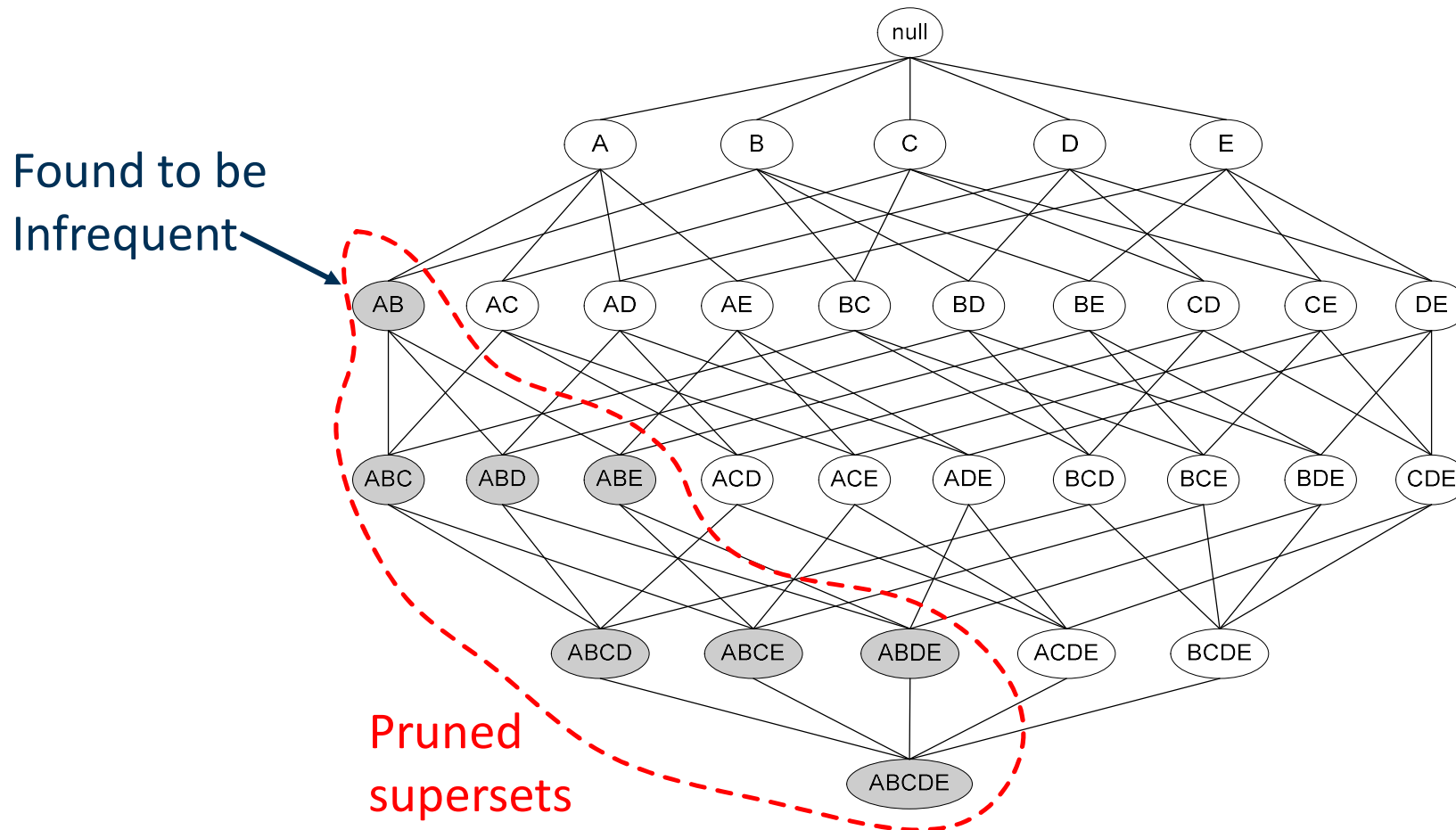
TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

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

- Consequence for frequent itemset search (aka Apriori principle):
  - If  $Y$  is frequent,  $X$  also has to be frequent
  - i.e.: **all subsets of frequent itemsets are frequent**

# Using the Apriori Principle for Pruning

- If an itemset is infrequent, then all of its supersets must also be infrequent



# The Apriori Algorithm

- Let  $k=1$
- Generate frequent itemsets of length 1
- Repeat until no new frequent itemsets are identified
  1. **Generate** length  $(k+1)$  candidate itemsets from length  $k$  frequent itemsets
  2. **Prune** candidate itemsets that can not be frequent because they contain subsets of length  $k$  that are infrequent (Apriori Principle)
  3. **Count** the support of each candidate by scanning the DB
  4. **Eliminate** candidates that are infrequent, leaving only those that are frequent

# Illustrating the Apriori Principle

Minimum Support Count = 3

## Items (1-itemsets)

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

No need to generate candidates involving Coke or Eggs



## Pairs (2-itemsets)

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

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Break, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

No need to generate candidate {Milk, Diaper, Beer} as count {Milk, Beer} = 2



## Triplets (3-itemsets)

Itemset	Count
{Bread, Milk, Diaper}	3

# Illustrating the Apriori Principle

- In the example, we had six items, and examined
  - Six 1-itemsets
  - Six 2-itemsets
  - One 3-itemset
  - i.e., **13 in total**
- vs. possible itemsets:  $2^6 = 64$  in total

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Break, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

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

Pruning strongly reduced the number of itemsets for which we need to scan the DB

# Example of Frequent Itemsets

Result History		FrequentItemSets				
Data	No. of Sets: 83	Size	Support ↓	Item 1	Item 2	Item 3
	Total Max. Size: 4	1	0.600	Asus EeePC		
Annotations	Min. Size: <input type="text" value="1"/>	1	0.500	LT Minimaus		
	Max. Size: <input type="text" value="4"/>	1	0.500	2 GB DDR3 RAM		
	Contains Item: <input type="text"/>	2	0.500	Asus EeePC	2 GB DDR3 RAM	
	<input type="button" value="Update View"/>	1	0.400	ThinkPad X220		
		1	0.400	Netbook-Schutzhülle		
		1	0.400	Lenovo Tablet Sleeve		
		1	0.400	LT Laser Maus		
		2	0.400	Asus EeePC	LT Minimaus	
		2	0.400	Asus EeePC	Netbook-Schutzhülle	
		2	0.400	2 GB DDR3 RAM	Netbook-Schutzhülle	
		3	0.400	Asus EeePC	2 GB DDR3 RAM	Netbook-Schutzhülle
		1	0.300	HP Laserjet P2055		
		1	0.300	HP CE50 Toner		
		2	0.300	LT Minimaus	2 GB DDR3 RAM	
		2	0.300	LT Minimaus	Netbook-Schutzhülle	
		2	0.300	ThinkPad X220	Lenovo Tablet Sleeve	
		2	0.300	HP Laserjet P2055	HP CE50 Toner	
		3	0.300	Asus EeePC	LT Minimaus	2 GB DDR3 RAM

# Application of Frequent Itemsets

1. Take top-k frequent itemsets of size 2 containing item A
2. Rank second item according to
  - profit made by selling item
  - whether you want to reduce number of items B in stock
  - knowledge about customer preferences
3. Offer special price for combination with top-ranked second item



Wird oft zusammen gekauft



Preis für beide: EUR 138,00

Beides in den Einkaufswagen

Verfügbarkeit und Versanddetails anzeigen

- ✓ **Dieser Artikel:** Introduction to Data Mining von Pang-Ning Tan Taschenbuch EUR 85,05
- ✓ Data Mining: Concepts and Techniques (Morgan Kaufmann Series in Data Management Systems)

# 3. From Frequent Itemsets to Rules

- Given a frequent itemset  $L$ , find all non-empty subsets  $f \subset L$  such that  $f \rightarrow L \setminus f$  satisfies the **minimum confidence** requirement

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Break, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

- Example Frequent Itemset  $L$ :
  - $\{Milk, Diaper, Beer\}$

- Example Rule:
  - $\underbrace{\{Milk, Diaper\}}_f \rightarrow \{Beer\}$

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

# Challenge: Combinatorial Explosion

- Given a 4-itemset  $\{A,B,C,D\}$ , we can generate
  - $\{A\} \rightarrow \{B,C,D\}$ ,  $\{A,B\} \rightarrow \{C,D\}$ ,  $\{B,C\} \rightarrow \{A,D\}$ ,  $\{C,D\} \rightarrow \{A,B\}$ ,  $\{A,B,C\} \rightarrow \{D\}$
  - $\{B\} \rightarrow \{A,C,D\}$ ,  $\{A,C\} \rightarrow \{B,D\}$ ,  $\{B,D\} \rightarrow \{A,C\}$ ,  $\{A,B,D\} \rightarrow \{C\}$
  - $\{C\} \rightarrow \{A,B,D\}$ ,  $\{A,D\} \rightarrow \{B,C\}$ ,  $\{A,C,D\} \rightarrow \{B\}$
  - $\{D\} \rightarrow \{A,B,C\}$ ,  $\{B,C,D\} \rightarrow \{A\}$
- i.e., a total of 14 rules for just one itemset!
- General number for a k-itemset:  $2^k - 2$ 
  - It's not  $2^k$  since we ignore  $\emptyset \rightarrow \{\dots\}$  and  $\{\dots\} \rightarrow \emptyset$

# Challenge: Combinatorial Explosion

- Wanted:  
another pruning trick like Apriori

- However

- $c(\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}) = 0.67$

- $c(\{\text{Milk}\} \rightarrow \{\text{Beer}\}) = 0.5$

- $c(\{\text{Diaper}\} \rightarrow \{\text{Beer}\}) = 0.8$

- $c(ABC \rightarrow D)$  can be larger or smaller than  $c(AB \rightarrow D)$

- In general, confidence does not have an anti-monotone property

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

# Challenge: Combinatorial Explosion

- But:  
confidence of rules generated  
from the **same itemset**  
has an anti-monotone property

- E.g.  $L = \{\text{Milk, Diaper, Beer}\}$ 
  - $\{\text{Milk, Diaper, Beer}\} \rightarrow \emptyset \ c=1.0$ 
    - $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\} \ c=0.67$
    - $\{\text{Milk}\} \rightarrow \{\text{Diaper, Beer}\} \ c=0.5$
    - $\{\text{Diaper}\} \rightarrow \{\text{Milk, Beer}\} \ c=0.5$
  - $\{\text{Milk, Beer}\} \rightarrow \{\text{Diaper}\} \ c=1.0$ 
    - $\{\text{Milk}\} \rightarrow \{\text{Diaper, Beer}\} \ c=0.5$
    - $\{\text{Beer}\} \rightarrow \{\text{Milk, Diaper}\} \ c=0.67$

- e.g.,  $L = \{A, B, C, D\}$ :  
$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Break, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

**Observation:** moving elements from antecedent to consequence (“left to right”) in the rule never increases confidence!

# Explanation

- Confidence is anti-monotone with respect to the number of items on the right-hand side (RHS) of the rule
  - i.e., “moving elements from left to right” cannot increase confidence
- Reason:

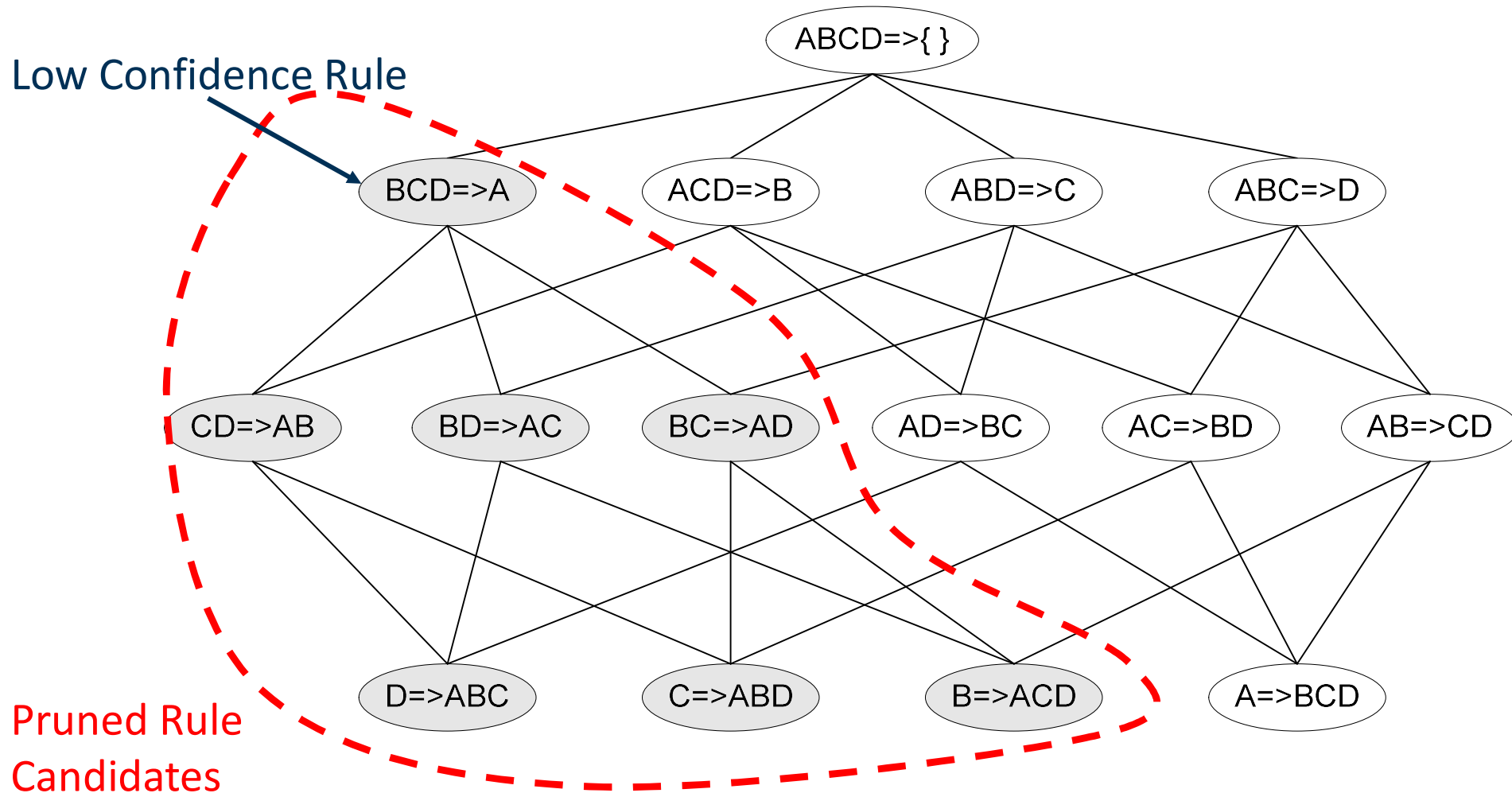
$$c(AB \rightarrow C) := \frac{s(ABC)}{s(AB)} \qquad c(A \rightarrow BC) := \frac{s(ABC)}{s(A)}$$

- Due to anti-monotone property of support, we know
$$s(AB) \leq s(A)$$
- Hence

$$c(AB \rightarrow C) \geq c(A \rightarrow BC)$$

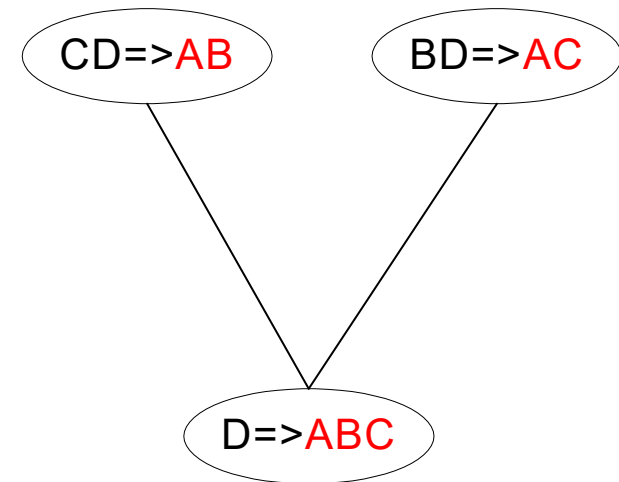
# Candidate Rule Pruning

- Moving elements from left to right cannot increase confidence



# Rule Generation for Apriori Algorithm

- Candidate rule is generated by merging two rules that share the same prefix in the rule consequent
  - $\text{join}(CD \rightarrow AB, BD \rightarrow AC)$  would produce the candidate rule  $D \rightarrow ABC$
  - Prune rule  $D \rightarrow ABC$  if one of its parent rules does not have high confidence (e.g.  $AD \rightarrow BC$ )



- All the required information for confidence computation has already been recorded during itemset generation
  - Thus, there is no need to see the data any more

$$c(X \rightarrow Y) = \frac{s(X \cup Y)}{s(X)}$$

## 4. Handling Continuous and Categorical Attributes

- How to apply association analysis to attributes that are not asymmetric binary variables?

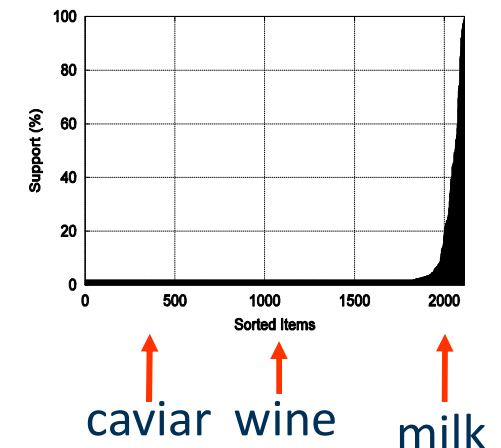
Session Id	Country	Session Length (sec)	Number of Web Pages viewed	Gender	Browser Type	Buy
1	USA	982	8	Male	Chrome	No
2	China	811	10	Female	Chrome	No
3	USA	2125	45	Female	Firefox	Yes
4	Germany	596	4	Male	IE	Yes
5	Australia	123	9	Male	Firefox	No
...	...	...	...	...	...	...

- Example Rule:

$\{\text{Number of Pages} \in [5,10) \wedge (\text{Browser}=\text{Firefox})\} \rightarrow \{\text{Buy} = \text{No}\}$

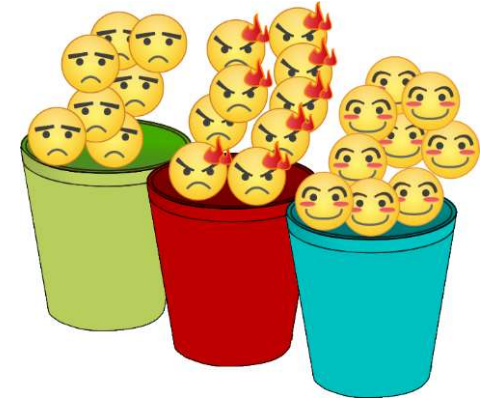
# Handling Categorical Attributes

- Transform categorical attribute into asymmetric binary variables
- Introduce a new “item” for each distinct attribute-value pair
  - e.g. replace “Browser Type” attribute with
    - attribute: “Browser Type = Chrome”
    - attribute: “Browser Type = Firefox”
    - .....
  - Python: OneHotEncoder
- Issues
  - What if attribute has many possible values?
    - Many of the attribute values may have very low support
    - Potential solution: aggregate low-support attribute values
  - What if distribution of attribute values is highly skewed?
    - Example: 95% of the visitors have Buy = No
    - Most of the items will be associated with (Buy=No) item
    - Potential solution: drop the highly frequent item



# Handling Continuous Attributes

- Transform continuous attribute into binary variables using discretization
  - Values of the attribute, e.g., age of a person:
    - 0, 4, 12, 16, 16, 18, 24, 26, 30
  - **Equal-interval binning** – for bin width of e.g., 10:
    - Bin 1: 0, 4                                     $[-\infty, 10)$  bin
    - Bin 2: 12, 16, 16, 18                     $[10, 20)$  bin
    - Bin 3: 24, 26, 30                          $[20, +\infty)$  bin
  - **Equal-frequency binning** – for bin density of e.g., 3:
    - Bin 1: 0, 4, 12                             $[-, 14)$  bin
    - Bin 2: 16, 16, 18                          $[14, 21)$  bin
    - Bin 3: 24, 26, 30                          $[21, +]$  bin



# Handling Continuous Attributes

- Issue: Size of the discretization intervals affects support & confidence
  - {Refund=No, (Income=\$51,251)} → {Cheat=No}
  - {Refund=No, (60K<= Income <=80K)} → {Cheat=No}
  - {Refund=No, (0K<= Income <=1B)} → {Cheat=No}
  
  - If intervals are too small
    - Itemsets may not have enough support
  - If intervals too large
    - Rules may not have enough confidence
    - e.g. combination of different age groups compared to a specific age group

# Association Analysis in Python

- Various packages exist
  - In the exercise, we'll use the Orange3 package
  - Frequent Itemset Generation

Python

```
from orangecontrib.associate.fpgrowth import frequent_itemsets  
  
# Calculate frequent itemsets  
itemsets = dict(frequent_itemsets(dataset.values, 0.20))
```

Min  
support

- Creating Association Rules

Python

```
from orangecontrib.associate.fpgrowth import association_rules  
  
# Calculate association rules from itemsets  
rules = association_rules(itemsets, 0.70)
```

Min  
confidence

## 5. Interestingness Measures

- Association rule algorithms tend to produce too many rules
  - Many of them are uninteresting or redundant
  - Redundant if  $\{A,B,C\} \rightarrow \{D\}$  and  $\{A,B\} \rightarrow \{D\}$  have same support & confidence
- Interestingness measures can be used to prune or rank the derived rules
- In the original formulation of association rules, support & confidence were the only interestingness measures used
- Later, various other measures have been proposed
  - We will have a look at one: Lift
  - See Tan/Steinbach/Kumar, Chapter 6.7

# Drawback of Confidence

Contingency table

	Coffee	$\overline{\text{Coffee}}$	
Tea	3	1	4
$\overline{\text{Tea}}$	15	1	16
	18	2	20



TID	Items
1	Coffee
2	Coffee
3	Coffee
4	Coffee
5	Coffee
6	Coffee
7	Coffee
8	Coffee
9	Coffee
10	Coffee
11	Coffee
12	Coffee
13	Coffee
14	Coffee
15	Coffee
16	Tea, Coffee
17	Tea, Coffee
18	Tea, Coffee
19	Tea
20	Bread

- Association Rule:  
Tea  $\rightarrow$  Coffee
- $\text{confidence}(\text{Tea} \rightarrow \text{Coffee}) = \frac{3}{4} = 0.75$
- **but**  $\text{support}(\text{Coffee}) = \frac{18}{20} = 0.9$
- Although confidence is high, rule is misleading as the fraction of coffee drinkers is higher than the confidence of the rule
  - We want  $\text{confidence}(X \rightarrow Y) > \text{support}(Y)$
  - otherwise rule is misleading as X reduces probability of Y

# Lift

- We discover a high confidence rule for tea  $\rightarrow$  coffee
  - 75% of all people who drink tea also drink coffee
  - Hypothesis: people who drink tea are likely to drink coffee
    - Implicitly: **more likely than all people**
- Test: Compare the confidence of the two rules
  - Rule: Tea  $\rightarrow$  coffee  $c(\text{tea} \rightarrow \text{coffee}) = \frac{s(\{\text{tea}\} \cup \{\text{coffee}\})}{s(\{\text{tea}\})}$
  - Default rule: all  $\rightarrow$  coffee  $c(\text{all} \rightarrow \text{coffee}) = \frac{s(\{\text{all}\} \cup \{\text{coffee}\})}{s(\{\text{all}\})} = \frac{s(\{\text{coffee}\})}{1} = s(\{\text{coffee}\})$
- We accept a rule iff its confidence is higher than the default rule

$$\frac{c(\text{tea} \rightarrow \text{coffee})}{c(\text{all} \rightarrow \text{coffee})} = \frac{c(\text{tea} \rightarrow \text{coffee})}{s(\{\text{coffee}\})} > 1$$

# Lift

- The lift of an association rule  $X \rightarrow Y$  is defined as:

$$c(X \rightarrow Y) = \frac{s(X \cup Y)}{s(X)}$$

$$Lift = \frac{c(X \rightarrow Y)}{s(Y)} = \frac{s(X \cup Y)}{s(X) * s(Y)}$$

- Confidence normalized by support of consequent
- Interpretation
  - if **lift** > **1**, then X and Y are positively correlated
  - if **lift** = **1**, then X and Y are independent
  - if **lift** < **1**, then X and Y are negatively correlated

# Lift Example

Contingency table

	Coffee	$\overline{\text{Coffee}}$	
Tea	3	1	4
$\overline{\text{Tea}}$	15	1	16
	18	2	20



TID	Items
1	Coffee
2	Coffee
3	Coffee
4	Coffee
5	Coffee
6	Coffee
7	Coffee
8	Coffee
9	Coffee
10	Coffee
11	Coffee
12	Coffee
13	Coffee
14	Coffee
15	Coffee
16	Tea, Coffee
17	Tea, Coffee
18	Tea, Coffee
19	Tea
20	Bread

- Association Rule:  
Tea  $\rightarrow$  Coffee

- confidence(Tea  $\rightarrow$  Coffee) =  $\frac{3}{4} = 0.75$

- **but** support(Coffee) =  $\frac{18}{20} = 0.9$

$$\begin{aligned} \text{Lift}(\text{Tea} \rightarrow \text{Coffee}) &= \frac{c(\text{tea} \rightarrow \text{coffee})}{c(\text{all} \rightarrow \text{coffee})} = \frac{c(\text{tea} \rightarrow \text{coffee})}{s(\{\text{coffee}\})} \\ &= \frac{0.75}{0.9} = 0.833 < 1 \end{aligned}$$

– lift < 1, therefore is negatively correlated and removed

# Conclusion

- The algorithm does the counting for you and finds patterns in the data
- You need to do the interpretation based on your knowledge about the application domain.
  - Which patterns are meaningful?
  - Which patterns are surprising?

# 6. Sequential Pattern Mining

- Association rule mining does not consider **the order of transactions**
- In many applications such orderings are relevant

Sequence Database	Sequence	Element (Transaction)	Event (Item)
Customer Data	Purchase history of a given customer	A set of items bought by a customer at time $t$	Books, diary products, CDs, etc
Web Server Logs	Browsing activity of a particular visitor	A collection of files viewed by a Web visitor after a single mouse click	Home page, index page, contact info, etc
Sensor Data	History of events generated by a given sensor	Events triggered by a sensor at time $t$	Types of alarms generated by sensors

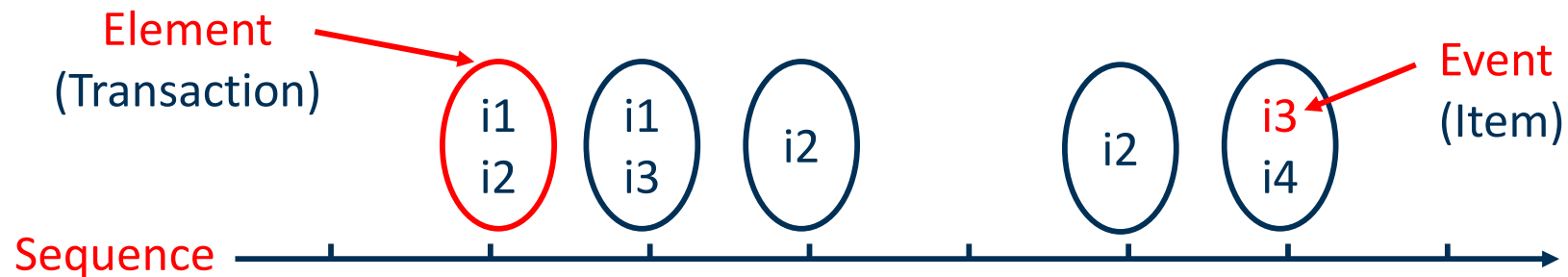
# Definition of a Sequence

- A **sequence** is an ordered list of elements (transactions)

$$s = \langle e_1 e_2 e_3 \dots \rangle$$

- Each **element** contains a collection of **events** (items)

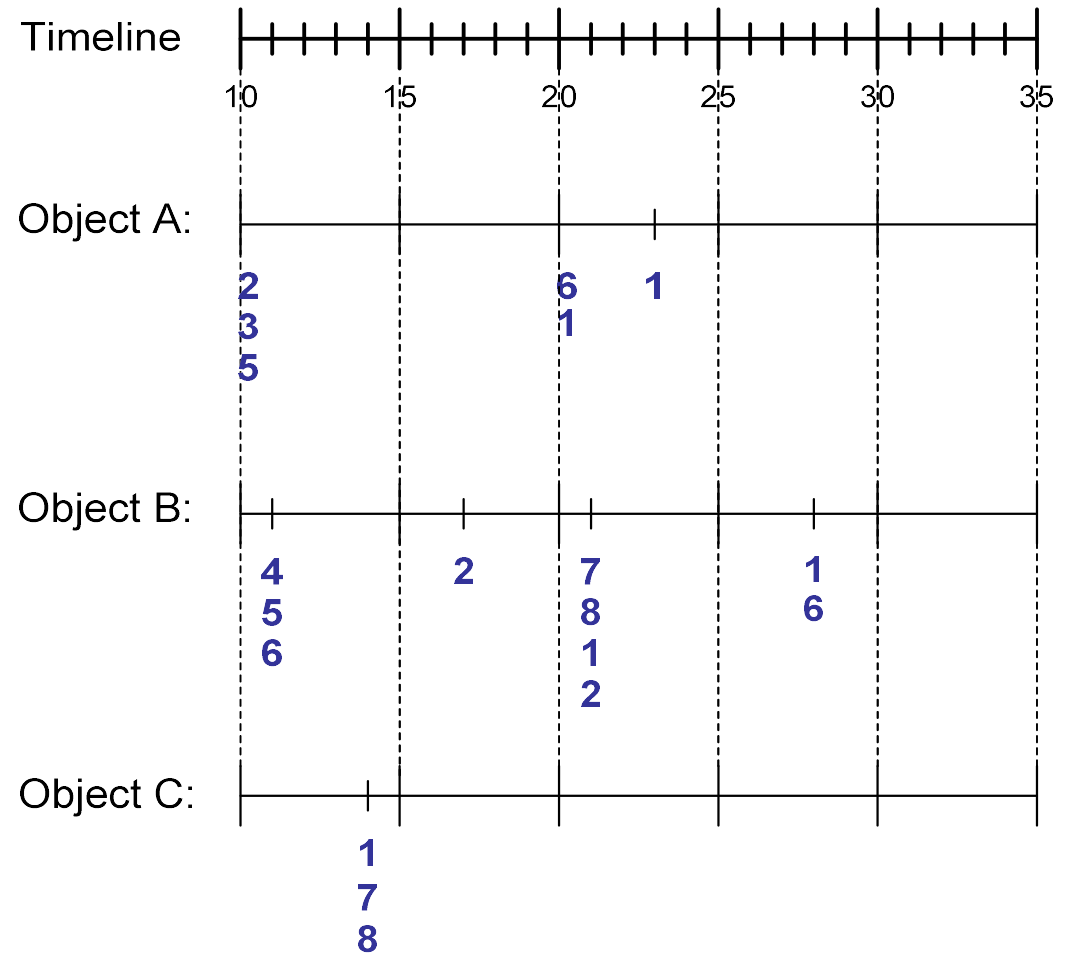
$$e_i = \{i_1, i_2, \dots, i_k\}$$



# Sequence Data

## Sequence Database:

Object	Timestamp	Events
A	10	2, 3, 5
A	20	6, 1
A	23	1
B	11	4, 5, 6
B	17	2
B	21	7, 8, 1, 2
B	28	1, 6
C	14	1, 8, 7



# Sequential Pattern Mining

- Formal Definition of a Subsequence
  - A sequence  $\langle a_1 a_2 \dots a_n \rangle$  is contained in another sequence  $\langle b_1 b_2 \dots b_m \rangle$  ( $m \geq n$ ) if there exist integers  $i_1 < i_2 < \dots < i_n$  such that  $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$

Data sequence $\langle b \rangle$	Subsequence $\langle a \rangle$	Contain?
$\langle \{2,4\} \{3,5,6\} \{8\} \rangle$	$\langle \{2\} \{3,5\} \rangle$	Yes
$\langle \{1,2\} \{3,4\} \rangle$	$\langle \{1\} \{2\} \rangle$	No
$\langle \{2,4\} \{2,4\} \{2,5\} \rangle$	$\langle \{2\} \{4\} \rangle$	Yes

- The **support** of a subsequence  $w$  is defined as the fraction of data sequences that contain  $w$
- A **sequential pattern** is a frequent subsequence (i.e., a subsequence whose support is  $\geq \text{minsup}$ )

# Sequential Pattern Mining

- Given:
  - a database of sequences
  - a user-specified minimum support threshold, *minsup*
- Task:
  - Find all subsequences with support  $\geq minsup$
- Challenge:
  - Very large number of candidate subsequences that need to be checked against the sequence database.
  - By applying the Apriori principle, the number of candidates can be pruned significantly.

# Determining the Candidate Subsequences

- Given  $n$  events:  $i_1, i_2, i_3, \dots, i_n$
- Candidate 1-subsequences:  
 $\langle \{i_1\} \rangle, \langle \{i_2\} \rangle, \langle \{i_3\} \rangle, \dots, \langle \{i_n\} \rangle$
- Candidate 2-subsequences:  
 $\langle \{i_1, i_2\} \rangle, \langle \{i_1, i_3\} \rangle, \dots, \langle \{i_1\} \{i_1\} \rangle, \langle \{i_1\} \{i_2\} \rangle, \dots, \langle \{i_{n-1}\} \{i_n\} \rangle$
- Candidate 3-subsequences:  
 $\langle \{i_1, i_2, i_3\} \rangle, \langle \{i_1, i_2, i_4\} \rangle, \dots, \langle \{i_1, i_2\} \{i_1\} \rangle, \langle \{i_1, i_2\} \{i_2\} \rangle, \dots,$   
 $\langle \{i_1\} \{i_1, i_2\} \rangle, \langle \{i_1\} \{i_1, i_3\} \rangle, \dots, \langle \{i_1\} \{i_1\} \{i_1\} \rangle, \langle \{i_1\} \{i_1\} \{i_2\} \rangle, \dots$

# Generalized Sequential Pattern Algorithm (GSP)

- Make the first pass over the sequence database  $D$  to yield all the 1-element frequent subsequences
- Repeat until no new frequent subsequences are found
  1. Candidate **Generation**:  
Merge pairs of frequent subsequences found in the  $(k-1)$ th pass to generate candidate sequences that contain  $k$  items
  2. Candidate **Pruning**:  
Prune candidate  $k$ -sequences that contain infrequent  $(k-1)$ -subsequences (**Apriori principle**)
  3. Support **Counting**:  
Make a new pass over the sequence database  $D$  to find the support for these candidate sequences
  4. Candidate **Elimination**:  
Eliminate candidate  $k$ -sequences whose actual support is less than *minsup*

# Candidate Generation

- Base case ( $k=2$ ):
  - Merging two frequent 1-sequences  $\langle\{i_1\}\rangle$  and  $\langle\{i_2\}\rangle$  will produce two candidate 2-sequences:  $\langle\{i_1\} \{i_2\}\rangle$  and  $\langle\{i_1 i_2\}\rangle$
- General case ( $k>2$ ):
  - A frequent  $(k-1)$ -sequence  $w_1$  is merged with another frequent  $(k-1)$ -sequence  $w_2$  to produce a candidate  $k$ -sequence if the subsequence obtained by **removing the first event in  $w_1$**  is the same as the subsequence obtained by **removing the last event in  $w_2$**
  - The resulting candidate after merging is given by the sequence  $w_1$  extended with the last event of  $w_2$ .
    - If the last two events in  $w_2$  belong to the same element, then the last event in  $w_2$  becomes part of the last element in  $w_1$
    - Otherwise, the last event in  $w_2$  becomes a separate element appended to the end of  $w_1$

# GSP Example

## Frequent 3-sequences

< {1} {2} {3} >  
< {1} {2 5} >  
< {1} {5} {3} >  
< {2} {3} {4} >  
< {2 5} {3} >  
< {3} {4} {5} >  
< {5} {3 4} >

## Candidate Generation

< {1} {2} {3} {4} >  
< {1} {2 5} {3} >  
< {1} {5} {3 4} >  
< {2} {3} {4} {5} >  
< {2 5} {3 4} >

## Candidate Pruning

< {1} {2 5} {3} >

- Only one 4-sequence survives the candidate pruning step
- All other 4-sequences are removed because they contain subsequences that are not part of the set of frequent 3-sequences.

# Timing Constraints

- Timing constraints allow us to pose **additional restrictions** on whether a sequence is counted to support a pattern or not.
- Motivating Example:
  - < {Statistics} {Database Systems} {Data Mining} >
  - < {Database Systems} {Statistics} {Data Mining} >
- We are interested in students that support the pattern  
< {Database Systems, Statistics} {Data Mining} >
- We don't care about the order of Database Systems and Statistics.
- We care about that the gap between these courses and Data Mining is not too long.

# Window Size

- Specifies a time window in the data sequence in which all events will be considered to belong to the same element.
- Given a candidate pattern:  $\langle \{a, c\} \rangle$
- Any data sequences that contain
  - $\langle \dots \{a\} \dots \{c\} \dots \rangle$ ,
  - $\langle \dots \{a\} \dots \{c\} \dots \rangle$  (where  $\text{time}(\{c\}) - \text{time}(\{a\}) \leq ws$ )
  - $\langle \dots \{c\} \dots \{a\} \dots \rangle$  (where  $\text{time}(\{a\}) - \text{time}(\{c\}) \leq ws$ )

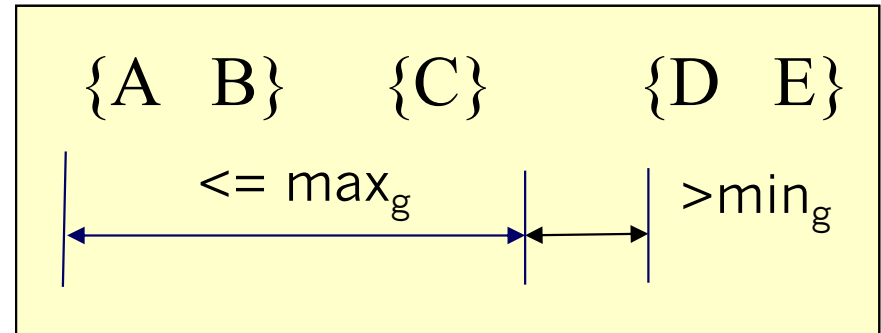
will contribute to the support count of the candidate pattern

- $\max_g = 2, \min_g = 0, ws = 1$

Data sequence	Subsequence	Contain?
$\langle \{2,4\} \{3,5,6\} \{4,7\} \{4,6\} \{8\} \rangle$	$\langle \{3\} \{5\} \rangle$	No
$\langle \{1\} \{2\} \{3\} \{4\} \{5\} \rangle$	$\langle \{1,2\} \{3\} \rangle$	Yes
$\langle \{1,2\} \{2,3\} \{3,4\} \{4,5\} \rangle$	$\langle \{1,2\} \{3,4\} \rangle$	Yes

# Max-Gap, Min-Gap

- Max-Gap: Sequence is counted if gap between consecutive elements is smaller than  $\max_g$
- Min-Gap: Sequence is counted if gap between consecutive elements is bigger than  $\min_g$
- $\max_g = 2, \min_g = 0$



Data sequence	Subsequence	Contain?
$\langle \{2,4\} \{3,5,6\} \{4,7\} \{4,5\} \{8\} \rangle$	$\langle \{6\} \{5\} \rangle$	Yes
$\langle \{1\} \{2\} \{3\} \{4\} \{5\} \rangle$	$\langle \{1\} \{4\} \rangle$	No
$\langle \{1\} \{2,3\} \{3,4\} \{4,5\} \rangle$	$\langle \{2\} \{3\} \{5\} \rangle$	Yes
$\langle \{1,2\} \{3\} \{2,3\} \{3,4\} \{2,4\} \{4,5\} \rangle$	$\langle \{1,2\} \{5\} \rangle$	No

# Literature for this Slideset

- Pang-Ning Tan, Michael Steinbach, Karpatne, Vipin Kumar:  
Introduction to Data Mining.  
2nd Edition. Pearson.
- Chapter 4: Association Analysis:  
Basic Concepts and  
Algorithms
- Chapter 7: Association Analysis:  
Advanced Concepts

