# Advanced Topics on Association Rules and Mining Sequence Data

Lecturer:  JERZY STEFANOWSKI

Institute of Computing Sciences

Poznan University of Technology

Poznan, Poland

Lectures 11

SE Master Course 2010

## Acknowledgments:

This lecture is based on the following resources - slides:
G.Piatetsky-Shapiro: Association Rules and Frequent Item Analysis.
and partly on two lectures
J.Han: Mining Association Rules in Large Databases;
Tan, Steinbach, Kumar: Introduction to Data Mining and my other notes.

# Outline

- Transactions

- Frequent itemsets

- Subset Property

- Association rules

- Applications

# Association rules

- Transaction data

- Market basket analysis



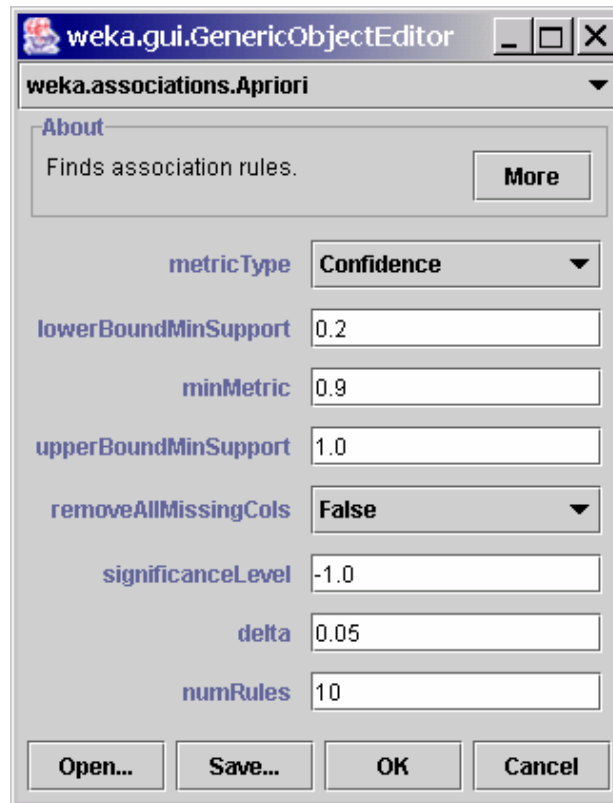| TID | Produce |
|-----|---------|
| 1 | MILK, BREAD, EGGS |
| 2 | BREAD, SUGAR |
| 3 | BREAD, CEREAL |
| 4 | MILK, BREAD, SUGAR |
| 5 | MILK, CEREAL |
| 6 | BREAD, CEREAL |
| 7 | MILK, CEREAL |
| 8 | MILK, BREAD, CEREAL, EGGS |
| 9 | MILK, BREAD, CEREAL |

- {Cereal, Milk} → Bread [sup=5%, conf=80%]

- Association rule:
  „80% of customers who buy *cereal* and *milk* also buy *bread* and 5% of customers buy all these products together"

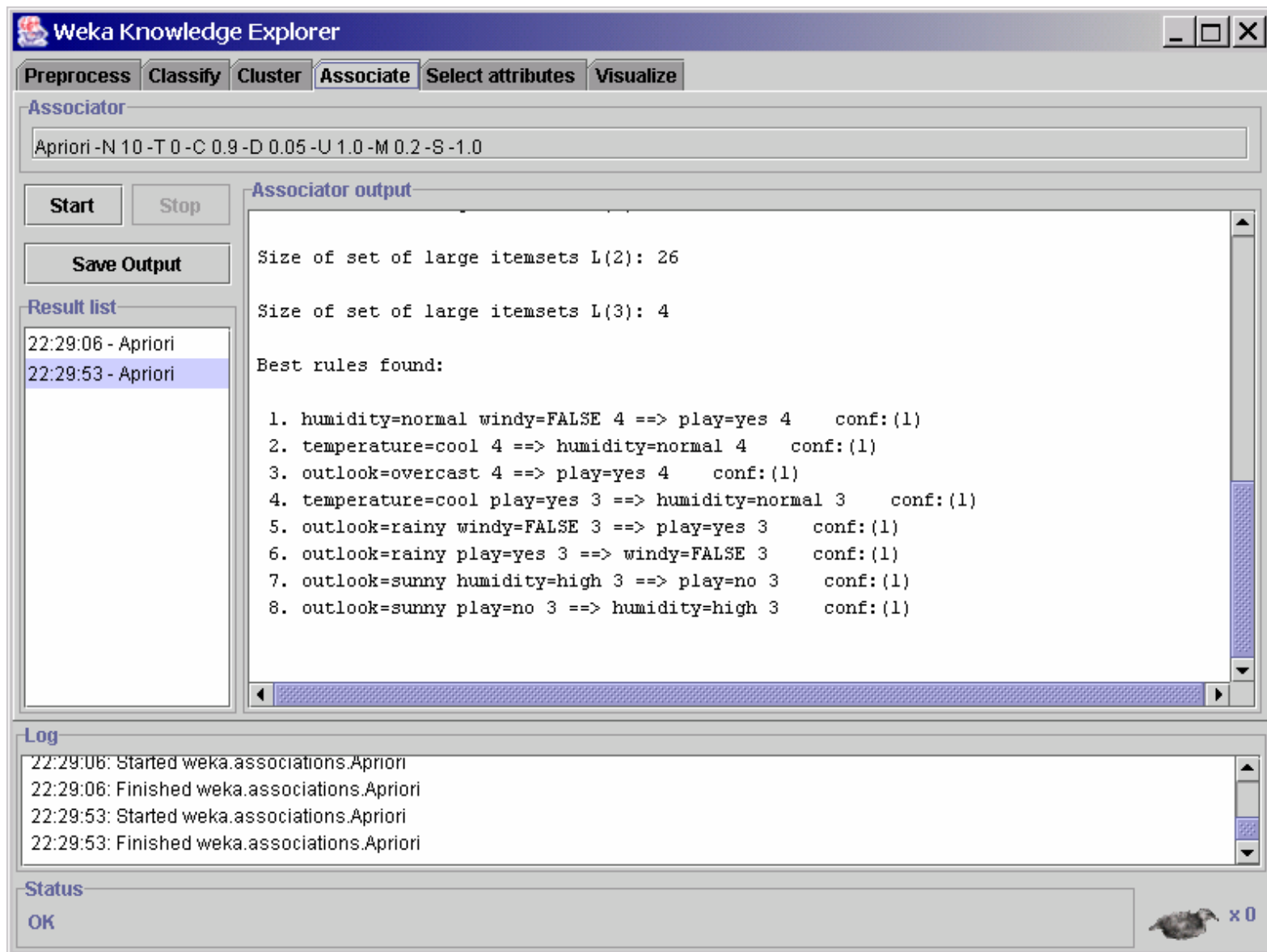Implication means co-occurrence, not causality!

# Weka associations

File: weather.nominal.arff
MinSupport: 0.2

# Weka associations: output



Weka Knowledge Explorer

Preprocess | Classify | Cluster | **Associate** | Select attributes | Visualize

**Associator**

Apriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.2 -S -1.0

Start | Stop

Save Output

**Result list**

22:29:06 - Apriori
22:29:53 - Apriori

**Associator output**

```
Size of set of large itemsets L(2): 26

Size of set of large itemsets L(3): 4

Best rules found:

1. humidity=normal windy=FALSE 4 ==> play=yes 4    conf:(1)
2. temperature=cool 4 ==> humidity=normal 4    conf:(1)
3. outlook=overcast 4 ==> play=yes 4    conf:(1)
4. temperature=cool play=yes 3 ==> humidity=normal 3    conf:(1)
5. outlook=rainy windy=FALSE 3 ==> play=yes 3    conf:(1)
6. outlook=rainy play=yes 3 ==> windy=FALSE 3    conf:(1)
7. outlook=sunny humidity=high 3 ==> play=no 3    conf:(1)
8. outlook=sunny play=no 3 ==> humidity=high 3    conf:(1)
```
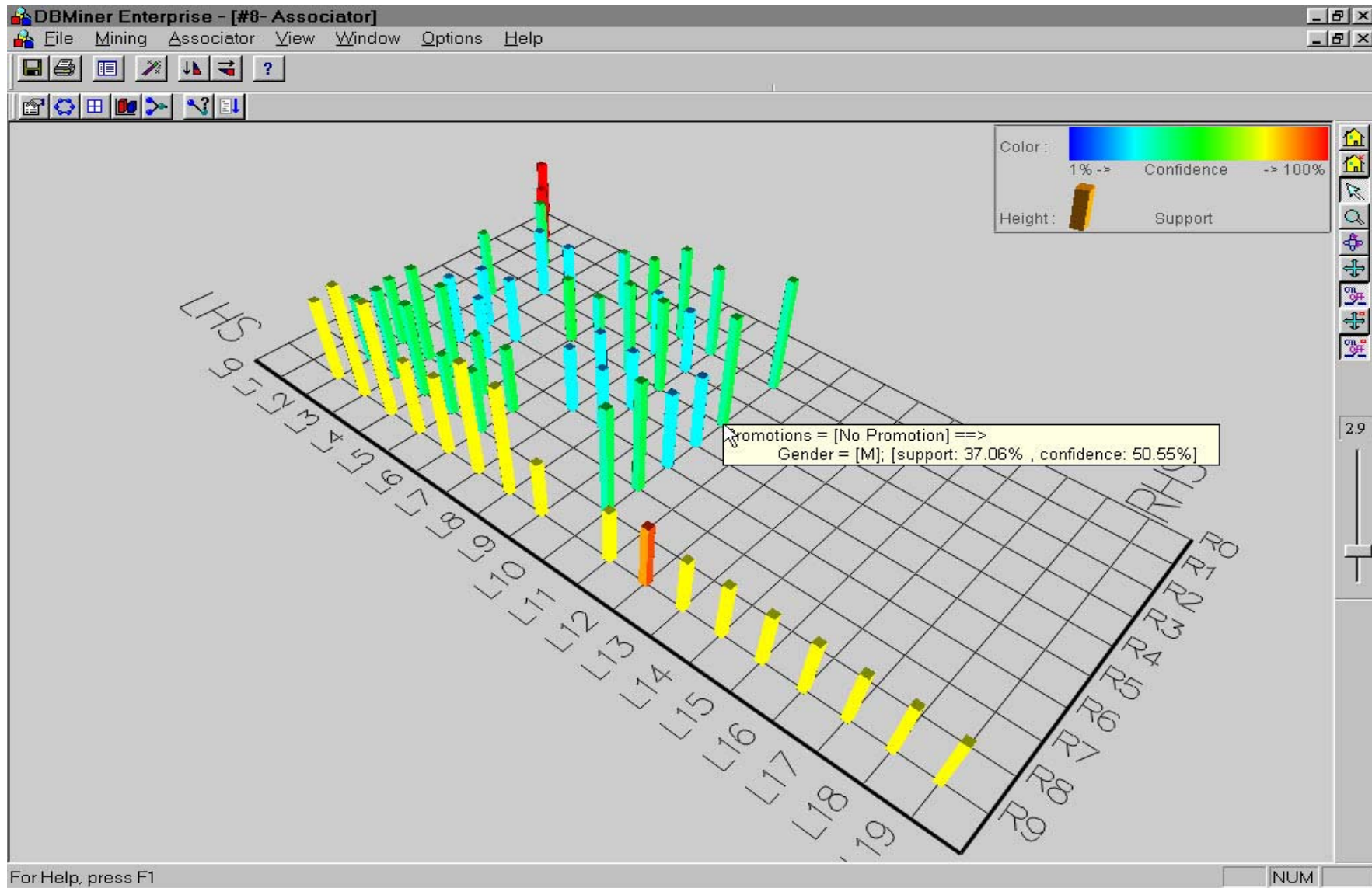
**Log**

22:29:06: Started weka.associations.Apriori
22:29:06: Finished weka.associations.Apriori
22:29:53: Started weka.associations.Apriori
22:29:53: Finished weka.associations.Apriori

**Status**

OK                                                                  x 0

# Presentation of Association Rules (Table Form )

| | Body | Implies | Head | Supp (%) | Conf (%) | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|
| 1 | cost(x) = '0.00~1000.00' | ==> | revenue(x) = '0.00~500.00' | 28.45 | 40.4 | | | | |
| 2 | cost(x) = '0.00~1000.00' | ==> | revenue(x) = '500.00~1000.00' | 20.46 | 29.05 | | | | |
| 3 | cost(x) = '0.00~1000.00' | ==> | order_qty(x) = '0.00~100.00' | 59.17 | 84.04 | | | | |
| 4 | cost(x) = '0.00~1000.00' | ==> | revenue(x) = '1000.00~1500.00' | 10.45 | 14.84 | | | | |
| 5 | cost(x) = '0.00~1000.00' | ==> | region(x) = 'United States' | 22.56 | 32.04 | | | | |
| 6 | cost(x) = '1000.00~2000.00' | ==> | order_qty(x) = '0.00~100.00' | 12.91 | 69.34 | | | | |
| 7 | order_qty(x) = '0.00~100.00' | ==> | revenue(x) = '0.00~500.00' | 28.45 | 34.54 | | | | |
| 8 | order_qty(x) = '0.00~100.00' | ==> | cost(x) = '1000.00~2000.00' | 12.91 | 15.67 | | | | |
| 9 | order_qty(x) = '0.00~100.00' | ==> | region(x) = 'United States' | 25.9 | 31.45 | | | | |
| 10 | order_qty(x) = '0.00~100.00' | ==> | cost(x) = '0.00~1000.00' | 59.17 | 71.86 | | | | |
| 11 | order_qty(x) = '0.00~100.00' | ==> | product_line(x) = 'Tents' | 13.52 | 16.42 | | | | |
| 12 | order_qty(x) = '0.00~100.00' | ==> | revenue(x) = '500.00~1000.00' | 19.67 | 23.88 | | | | |
| 13 | product_line(x) = 'Tents' | ==> | order_qty(x) = '0.00~100.00' | 13.52 | 98.72 | | | | |
| 14 | region(x) = 'United States' | ==> | order_qty(x) = '0.00~100.00' | 25.9 | 81.94 | | | | |
| 15 | region(x) = 'United States' | ==> | cost(x) = '0.00~1000.00' | 22.56 | 71.39 | | | | |
| 16 | revenue(x) = '0.00~500.00' | ==> | cost(x) = '0.00~1000.00' | 28.45 | 100 | | | | |
| 17 | revenue(x) = '0.00~500.00' | ==> | order_qty(x) = '0.00~100.00' | 28.45 | 100 | | | | |
| 18 | revenue(x) = '1000.00~1500.00' | ==> | cost(x) = '0.00~1000.00' | 10.45 | 96.75 | | | | |
| 19 | revenue(x) = '500.00~1000.00' | ==> | cost(x) = '0.00~1000.00' | 20.46 | 100 | | | | |
| 20 | revenue(x) = '500.00~1000.00' | ==> | order_qty(x) = '0.00~100.00' | 19.67 | 96.14 | | | | |
| 21 | | | | | | | | | |
| 22 | | | | | | | | | |
| 23 | cost(x) = '0.00~1000.00' | ==> | revenue(x) = '0.00~500.00' AND order_qty(x) = '0.00~100.00' | 28.45 | 40.4 | | | | |
| 24 | cost(x) = '0.00~1000.00' | ==> | revenue(x) = '0.00~500.00' AND order_qty(x) = '0.00~100.00' | 28.45 | 40.4 | | | | |
| 25 | cost(x) = '0.00~1000.00' | ==> | revenue(x) = '500.00~1000.00' AND order_qty(x) = '0.00~100.00' | 19.67 | 27.93 | | | | |
| 26 | cost(x) = '0.00~1000.00' | ==> | revenue(x) = '500.00~1000.00' AND order_qty(x) = '0.00~100.00' | 19.67 | 27.93 | | | | |
| 27 | cost(x) = '0.00~1000.00' AND order_qty(x) = '0.00~100.00' | ==> | revenue(x) = '500.00~1000.00' | 19.67 | 33.23 | | | | |

Sheet1

# Visualization of Association Rules: Plane Graph

# Filtering Association Rules

- Finding Association Rules is just the beginning in a datamining effort.

- Problem: any large dataset can lead to a very large number of association rules, even with reasonable Min Confidence and Support

  - Many of these rules are uninteresting, trivial or redundant

- Trivial rule example:

  - pregnant $\rightarrow$ female with accuracy 1!

- Challenge is to select potentially interesting rules

- Finding Association rules is a kind of Exploratory Data Analysis

# Need for interestingness measures

- In the original formulation of association rules, support & confidence are the only measures used

- Confidence by itself is not sufficient
  - e.g. if all transactions include Z, then
  - any rule I => Z will have confidence 100%.

- Other interestingness measures are necessary to filter rules!

# Computing Interestingness Measure

- Given a rule $X \rightarrow Y$, information needed to compute rule interestingness can be obtained from a contingency table

Contingency table for $X \rightarrow Y$

|  | Y | $\overline{Y}$ |  |
|---|---|---|---|
| X | $f_{11}$ | $f_{10}$ | $f_{1+}$ |
| $\overline{X}$ | $f_{01}$ | $f_{00}$ | $f_{o+}$ |
|  | $f_{+1}$ | $f_{+0}$ | $|T|$ |

$f_{11}$: support of X and Y
$f_{10}$: support of X and $\overline{Y}$
$f_{01}$: support of $\overline{X}$ and Y
$f_{00}$: support of $\overline{X}$ and $\overline{Y}$

Used to define various measures

◆ support, confidence, lift, Gini, Piatetsky, J-measure, etc.

# Interestingness Measure: Correlations and Lift

- *play basketball* $\Rightarrow$ *eat cereal* [40%, 66.7%] is misleading

  - The overall percentage of students eating cereal is 75% which is higher than 66.7%.

- *play basketball* $\Rightarrow$ *not eat cereal* [20%, 33.3%] is more accurate, although with lower support and confidence

- Measure of dependent/correlated events: lift or corr, …

$$corr_{A,B} = \frac{P(A \cup B)}{P(A)P(B)}$$

| | Basketball | Not basketball | Sum (row) |
|---|---|---|---|
| Cereal | 2000 | 1750 | 3750 |
| Not cereal | 1000 | 250 | 1250 |
| Sum(col.) | 3000 | 2000 | 5000 |

# Statistical Independence

- Population of 1000 students

  - 600 students know how to swim (S)

  - 700 students know how to bike (B)

  - 420 students know how to swim and bike (S,B)

  - $P(S \land B) = 420/1000 = 0.42$

  - $P(S) \times P(B) = 0.6 \times 0.7 = 0.42$

  - $P(S \land B) = P(S) \times P(B)$ => Statistical independence

  - $P(S \land B) > P(S) \times P(B)$ => Positively correlated

  - $P(S \land B) < P(S) \times P(B)$ => Negatively correlated

# Association Rule LIFT

- The *lift* of an association rule $I => J$ is defined as:

  - lift = P(J|I) / P(J)

  - Note, P(J) = (support of J) / (no. of transactions)

  - ratio of confidence to expected confidence


- Interpretation:

  - if  lift > 1, then I and J are positively correlated

    lift < 1, then I are J are negatively correlated.

    lift = 1, then I and J are independent.

# Illustrative Example

| | Coffee | $\overline{\text{Coffee}}$ | |
|---|---|---|---|
| Tea | 15 | 5 | 20 |
| $\overline{\text{Tea}}$ | 75 | 5 | 80 |
| | 90 | 10 | 100 |

Drawback of using confidence only!

## Association Rule: Tea → Coffee

Confidence= P(Coffee|Tea) = 0.75

but P(Coffee) = 0.9

⇒ Although confidence is high, rule is misleading

⇒ P(Coffee|Tea) = 0.9375

# Example: Lift/Interest

| | 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) = 0.75

but P(Coffee) = 0.9

$\Rightarrow$ Lift = 0.75/0.9= 0.8333 (< 1, therefore is negatively associated)

# Statistical-based Measures

- Measures that take into account statistical dependence

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

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

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

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

# Drawback of Lift & Interest

$$X \rightarrow Y$$

|      | 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 |

$$P(X \cap Y) = 10/100 = P(X) = P(Y)$$

$$Lift = \frac{0.1}{(0.1)(0.1)} = 10 \qquad\qquad Lift = \frac{0.9}{(0.9)(0.9)} = 1.11$$

**Statistical independence:**

**If P(X,Y)=P(X)P(Y)  => Lift = 1**

# 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**

There are lots of measures proposed in the literature

Some measures are good for certain applications, but not for others

What criteria should we use to determine whether a measure is good or bad?

What about Apriori-style support based pruning? How does it affect these measures?

| # | 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 $(J)$ | $\max\left(P(A,B)\log(\frac{P(B\mid A)}{P(B)})+P(A\overline{B})\log(\frac{P(\overline{B}\mid A)}{P(\overline{B})}),\right.$ $\left. P(A,B)\log(\frac{P(A\mid B)}{P(A)})+P(\overline{A}B)\log(\frac{P(\overline{A}\mid B)}{P(A)})\right)$ |
| 9 | Gini index $(G)$ | $\max\left(P(A)[P(B\mid A)^2+P(\overline{B}\mid A)^2]+P(\overline{A})[P(B\mid\overline{A})^2+P(\overline{B}\mid\overline{A})^2]\right.$ $-P(B)^2-P(\overline{B})^2,$ $P(B)[P(A\mid B)^2+P(\overline{A}\mid B)^2]+P(\overline{B})[P(A\mid\overline{B})^2+P(\overline{A}\mid\overline{B})^2]$ $\left.-P(A)^2-P(\overline{A})^2\right)$ |
| 10 | Support $(s)$ | $P(A,B)$ |
| 11 | Confidence $(c)$ | $\max(P(B\mid A),P(A\mid B))$ |
| 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 $(F)$ | $\max\left(\frac{P(B\mid A)-P(B)}{1-P(B)},\frac{P(A\mid B)-P(A)}{1-P(A)}\right)$ |
| 18 | Added Value $(AV)$ | $\max(P(B\mid A)-P(B),P(A\mid B)-P(A))$ |
| 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 $(K)$ | $\sqrt{P(A,B)}\max(P(B\mid A)-P(B),P(A\mid B)-P(A))$ |

# Properties of A Good Measure

- **Piatetsky-Shapiro**:
  3 properties a good measure M must satisfy:

  - $M(A,B) = 0$ if A and B are statistically independent

  - $M(A,B)$ increase monotonically with $P(A,B)$ when $P(A)$ and $P(B)$ remain unchanged

  - $M(A,B)$ decreases monotonically with $P(A)$ [or $P(B)$] when $P(A,B)$ and $P(B)$ [or $P(A)$] remain unchanged

# Alternative approaches

- Multiple criteria approaches to many evaluation measures (Pareto border of the set of rules)

- Specific systems based on interaction with advanced users – directing the search

  - Templates as to the syntax

  - Other specifications for rules

# Manila, Toivonen Finding Interesting Association Rules



Figure 1: Rule Visualizer / Rule Selection.

# Visualization of rules
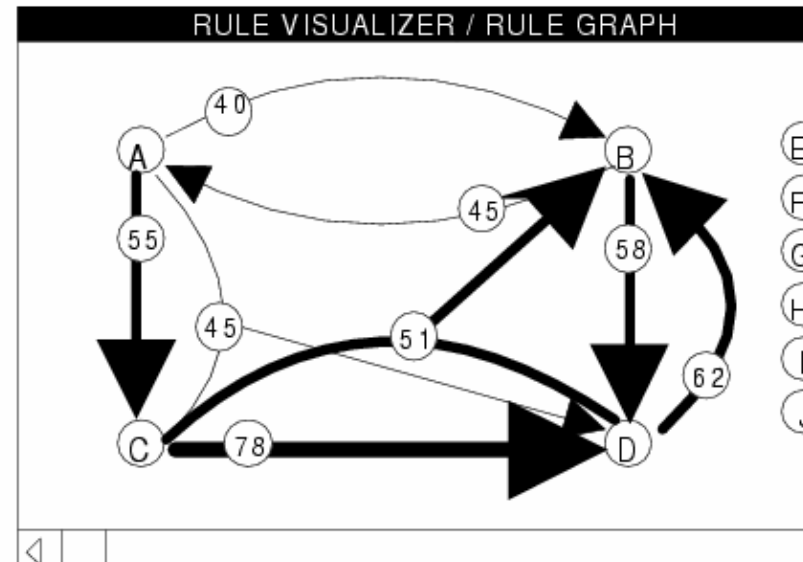


Figure 3: Rule Visualizer / Rule Browsing.
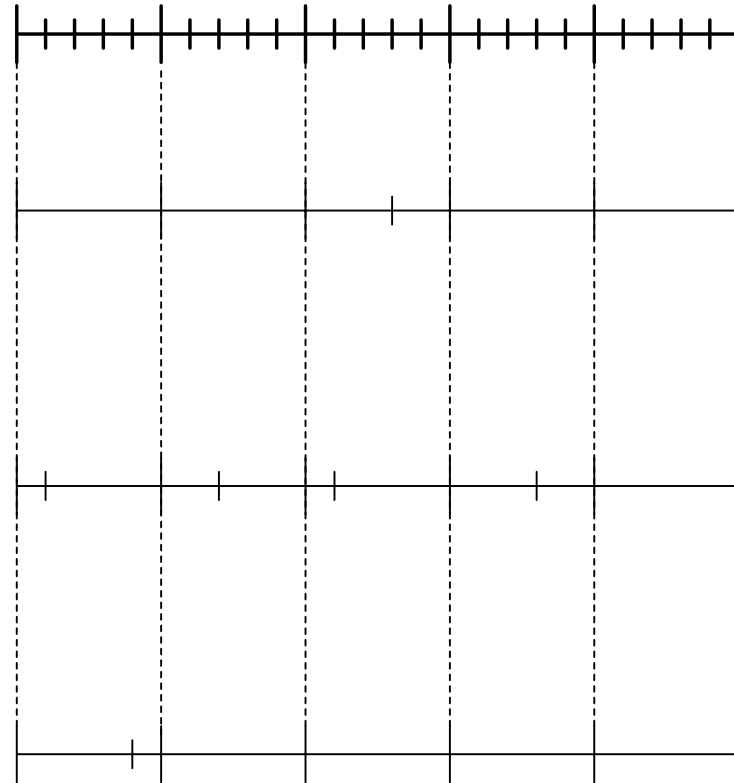


Figure 4: Rule Visualizer / Rule Graph.

# Mining sequence data

Another important problem strongly inspired by frequent itemsets and association rules!

# 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 |

# Sequence Databases and Sequential Pattern Analysis

- Transaction databases, time-series databases vs. sequence databases

- Frequent patterns vs. (frequent) sequential patterns

- Applications of sequential pattern mining

  - Customer shopping sequences:

    - First buy computer, then CD-ROM, and then digital camera, within 3 months.

  - Medical treatment, natural disasters (e.g., earthquakes), science & engineering processes, stocks and markets, etc.

  - Telephone calling patterns, Weblog click streams

  - DNA sequences and gene structures

# Examples of Sequence Data

| Sequence Database | Sequence | Element (Transaction) | Event (Item) |
|---|---|---|---|
| Customer | Purchase history of a given customer | A set of items bought by a customer at time t | Books, diary products, CDs, etc |
| Web Data | Browsing activity of a particular Web visitor | A collection of files viewed by a Web visitor after a single mouse click | Home page, index page, contact info, etc |
| Event data | History of events generated by a given sensor | Events triggered by a sensor at time t | Types of alarms generated by sensors |
| Genome sequences | DNA sequence of a particular species | An element of the DNA sequence | Bases A,T,G,C |

Element
(Transaction)

Event
(Item)

Sequence

E1 E2    E1 E3    E2         E2    E3 E4

# Formal Definition of a Sequence

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

$$s = < e_1\ e_2\ e_3\ ... >$$

  - Each element contains a collection of events (items)

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

  - Each element is attributed to a specific time or location

- Length of a sequence, $|s|$, is given by the number of elements of the sequence

- A k-sequence is a sequence that contains k events (items)

# Examples of Sequence

- Web sequence:

  < {Homepage}  {Electronics}  {Digital Cameras}  {Canon Digital Camera} {Shopping Cart}  {Order Confirmation}  {Return to Shopping} >

- Sequence of initiating events causing the nuclear accident at 3-mile Island:
  (http://stellar-one.com/nuclear/staff_reports/summary_SOE_the_initiating_event.htm)

  <   {clogged resin} {outlet valve closure} {loss of feedwater} {condenser polisher outlet valve shut} {booster pumps trip} {main waterpump trips} {main turbine trips} {reactor pressure increases}>

- Sequence of books checked out at a library:

  <{Fellowship of the Ring} {The Two Towers}  {Return of the King}>

# What Is Sequential Pattern Mining?

- Given a set of sequences, find the complete set of *frequent* subsequences

A *sequence* : < (ef) (ab) (df) c b >

A *sequence database*

| SID | sequence |
|-----|----------|
| 10 | <a(abc)(ac)d(cf)> |
| 20 | <(ad)c(bc)(ae)> |
| 30 | <(ef)(ab)(df)cb> |
| 40 | <eg(af)cbc> |

An element may contain a set of items. Items within an element are unordered and we list them alphabetically.

<a(bc)dc> is a *subsequence* of <a(abc)(ac)d(cf)>

Given *support threshold* *min_sup* =2, <(ab)c> is a *sequential pattern*
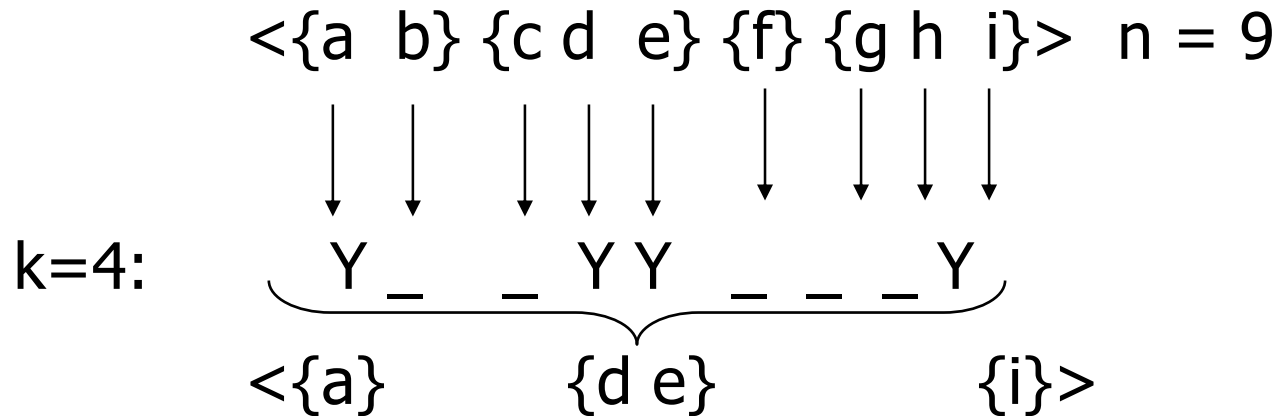
# Sequential Pattern Mining: Definition

- Given:

    - a database of sequences

    - a user-specified minimum support threshold, *minsup*


- Task:

    - Find all subsequences with support ≥ *minsup*

# Sequential Pattern Mining: Challenge

- Given a sequence:   <{a b} {c d e} {f} {g h i}>

  - Examples of subsequences:

    <{a} {c d} {f} {g} >, < {c d e} >, < {b} {g} >, etc.

- How many k-subsequences can be extracted from a given n-sequence?

$$<\{a\ b\}\ \{c\ d\ e\}\ \{f\}\ \{g\ h\ i\}>\ \ n = 9$$

k=4:     Y _   _ Y Y   _ _ _ Y

<{a}        {d e}        {i}>

Answer :

$$\binom{n}{k} = \binom{9}{4} = 126$$

# Challenges on Sequential Pattern Mining

- A huge number of possible sequential patterns are hidden in databases

- A mining algorithm should

  - find the complete set of patterns, when possible, satisfying the minimum support (frequency) threshold

  - be highly efficient, scalable, involving only a small number of database scans

  - be able to incorporate various kinds of user-specific constraints

# Studies on Sequential Pattern Mining

- Concept introduction and an initial Apriori-like algorithm
    - R. Agrawal & R. Srikant. "Mining sequential patterns," ICDE'95
- GSP—An Apriori-based, influential mining method (developed at IBM Almaden)
    - R. Srikant & R. Agrawal. "Mining sequential patterns: Generalizations and performance improvements," EDBT'96
- FreeSpan and PrefixSpan (Han et al.@KDD'00; Pei, et al.@ICDE'01)
    - Projection-based
    - But only prefix-based projection: less projections and quickly shrinking sequences
- Vertical format-based mining: SPADE (Zaki00)

# A Basic Property of Sequential Patterns: Apriori like approach

- A basic property: Apriori (Agrawal & Sirkant'94)

  - If a sequence S is not frequent

  - Then, none of the super-sequences of S is frequent

  - E.g, <hb> is infrequent → so do <hab> and <(ah)b>

| Seq. ID | Sequence |
|---------|----------|
| 10 | <(bd)cb(ac)> |
| 20 | <(bf)(ce)b(fg)> |
| 30 | <(ah)(bf)abf> |
| 40 | <(be)(ce)d> |
| 50 | <a(bd)bcb(ade)> |

Given *support threshold*
*min_sup* =2

# GSP—A Generalized Sequential Pattern Mining Algorithm

- GSP (Generalized Sequential Pattern) mining algorithm

    - proposed by Agrawal and Srikant, EDBT'96

- Outline of the method

    - Initially, every item in DB is a candidate of length-1

    - for each level (i.e., sequences of length-k) do

        - scan database to collect support count for each candidate sequence

        - generate candidate length-(k+1) sequences from length-k frequent sequences using Apriori

    - repeat until no frequent sequence or no candidate can be found

- Major strength: Candidate pruning by Apriori

# Performance on Data Set Gazelle

# Multidimesional sequentianl patterns

- Sequential patterns are useful

  - "free internet access $\rightarrow$ buy package 1 $\rightarrow$ upgrade to package 2"

  - Marketing, product design & development

- Problems: lack of focus

  - Various groups of customers may have different patterns

- MD-sequential pattern mining: integrate multi-dimensional analysis and sequential pattern mining

# An example of Multidim. Contxt sequential pattern

**Sequence /customer context**:

*Monthly earnings, Martial status, Profession, Age*

**Transaction context**:

*Time from money supply,*

*Day of the weak when action done*

**User actions**:

SD –receive money, TM – transfer WM – withdraw money, CD – create time deposit, RD – cancel this deposit

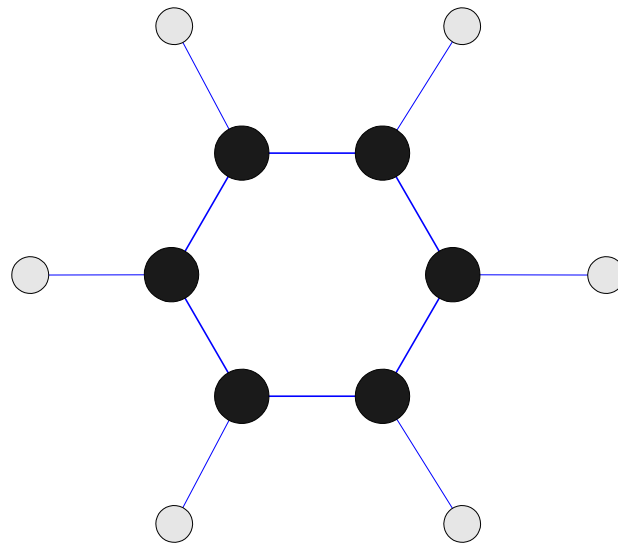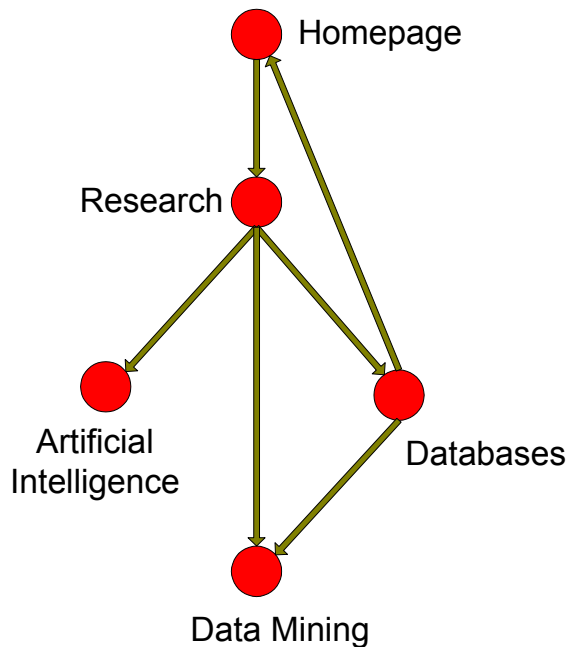| Sequences: | | |
|---|---|---|
| SID1 | (2,Friday) | {TM,CD} |
| (4200,married,tech,24) | (4,Sunday) | {WM} |
| | (20,Saturday) | {RD,WM,TM} |
| SID2 | (3,Tuesday) | {TM,CD,WM} |
| (4000,married,tech,22) | (7,Sunday) | {WM,CD} |
| | (20,Saturday) | {RD,WM} |
| | (1,Tuesday) | {TM,CD} |
| SID3 | (3,Monday) | {CD,TM,WM} |
| (1500,single,retired,70) | (10,Monday) | {CD,TM,WM} |
| | (16,Sunday) | {WM} |
| ... | | ... |

Examples of patterns:

- Traditional sequential pattern:

   <{TM,CD},{WM},{WM,RD}>

- **Extended context sequential pattern:**

   (4000,married,*,*)<(3,*){TM,CD},(*,Sunday){WM},(20,*){WM,RD}>

# Frequent Subgraph Mining

- Extend association rule mining to finding frequent subgraphs

- Useful for Web Mining, computational chemistry, bioinformatics, spatial data sets, etc

# Applications

- Market basket analysis
  - Store layout, client offers
- This analysis is applicable whenever a customer purchases multiple things in proximity
  - telecommunication (each customer is a transaction containing the set of phone calls)
  - weather analysis (each time interval is a transaction containing the set of observed events)
  - credit cards
  - banking services
  - medical treatments
- Finding unusual events
  - WSARE – What is Strange About Recent Events
- …

# Conclusions

- Association rule mining

    - probably the most significant contribution from the database community in KDD

    - A large number of papers have been published

- Many interesting issues have been explored

- An interesting research direction

    - Association analysis in other types of data: sequence data, spatial data, multimedia data, time series data, etc.

# Summary

- Frequent itemsets

- Association rules

- Subset property

- Apriori algorithm

- Extensions of this algorithm

- Evaluation of association rules

- Sequence patterns

# Any questions, remarks?