

# COMP3420: Advanced Databases and Data Mining

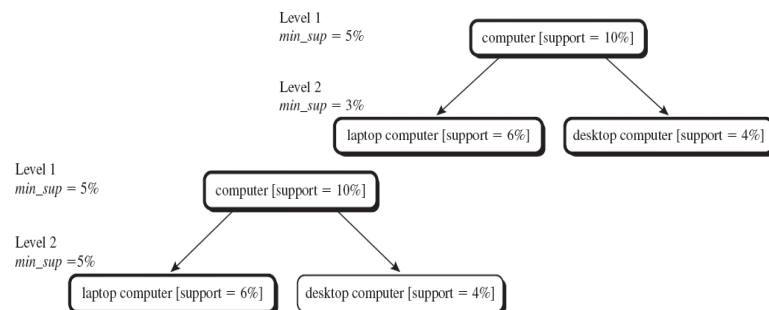
## Advanced association mining

## Lecture outline

- Mining various kinds of association rules
  - Multi-level association
  - Multi-dimensional association
  - Quantitative association
  - Interesting correlation patterns
- Constraints based mining
- Interestingness measure: Correlation (Lift)
  - More interestingness measures
- Visualisation of association rules

## Multi-level association mining

- Items often form hierarchies
- Items at lower levels are expected to have lower support
  - Flexible *support* setting (uniform, reduced, or group-based (user specific))



Source: Han and Kamber, DM Book, 2<sup>nd</sup> Ed. (Copyright © 2006 Elsevier Inc.)

## Multi-level association mining (2)

- Some rules may be redundant due to *ancestor* relationships between items
- For example:
  - $buys(X, \text{'milk'}) \Rightarrow buys(X, \text{'bread'})$  [ $s=8\%$ ,  $c=70\%$ ]
  - $buys(X, \text{'skim milk'}) \Rightarrow buys(X, \text{'bread'})$  [ $s=2\%$ ,  $c=72\%$ ]
  - The first rule is said to be an *ancestor* of the second rule
- A rule is redundant if its support is close to the “expected” value, based on the rule’s ancestor
  - For example, if around 25% of all milk purchased is ‘skim milk’, then the second rule above is redundant, as it has a  $\frac{1}{4}$  of the support of the first, more general rule (and similar confidence)

## Multi-dimensional association mining

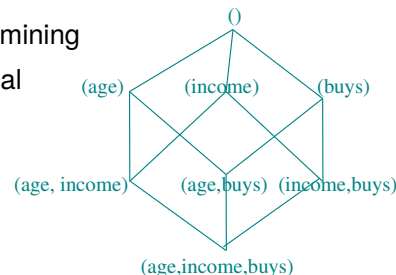
- Single-dimensional rules:  $buys(X, 'milk') \Rightarrow buys(X, 'bread')$
- Multi-dimensional rules: Two or more dimensions or predicates (or attributes)
  - Inter-dimension association rules (*no repeated predicates*):  
 $age(X, '19-25')$  and  $occupation(X, 'student') \Rightarrow buys(X, 'coke')$
  - Hybrid-dimension association rules (*repeated predicates*):  
 $age(X, '19-25')$  and  $buys(X, 'popcorn') \Rightarrow buys(X, 'coke')$
- Categorical Attributes: finite number of possible values, no ordering among values (data cube approach)
- Quantitative Attributes: numeric, implicit ordering among values (discretisation, clustering, etc.)

## Quantitative association mining

- Techniques can be categorised by how numerical attributes, such as *age* or *income*, are treated
- Static discretisation based on predefined concept hierarchies (data cube methods)
- Dynamic discretisation based on data distribution
  - $A_{quant1}$  and  $A_{quant2} \Rightarrow A_{cat}$
  - Example:  $age(X, '19-25')$  and  $income(X, '40K-60K') \Rightarrow buys(X, 'HDTV')$
- For quantitative rules, do discretisation such that (for example) the confidence of the rules mined is maximised

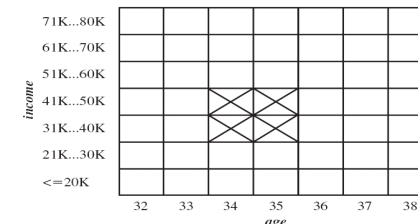
## Static discretisation of quantitative attributes

- Discretised prior to mining using concept hierarchy
- Numeric values are replaced by ranges
- In relational database, finding all frequent  $k$ -item sets will require  $k$  database scans
- Data cube is well suited for mining
- The cells of an  $n$ -dimensional cuboid correspond to the item or *predicate sets*
- Mining from data cubes can be much faster



## Dynamic discretisation of quantitative attributes

- Mapping of pairs of quantitative attributes into a 2-dimensional grid, such that categorical attribute conditions are satisfied
- The grid is then searched for clusters of points from which association rules are generated
- For example:  
 $age(X, '34-35')$  and  $income(X, '31K-50K') \Rightarrow buys(X, 'HDTV')$



## Mining interesting correlation patterns

- Flexible support
  - Some items might be very rare but are valuable (like diamonds)
  - Customise  $support_{min}$  specification and application
- Top- $k$  frequent patterns
  - It can be hard to specify  $support_{min}$ , but top- $k$  rules with  $length_{min}$  are more desirable
  - Achievable using special data structures, like Frequent-Pattern (FP) tree
  - Dynamically raise  $support_{min}$  during FP-tree construction phase, and select most promising to mine
- Many different algorithms and data structures have been developed to allow efficient mining of associations and rules (more in the text book)

## Constraint based mining

- Finding *all* the frequent rules or patterns in a database autonomously is unrealistic
  - The rules / patterns could be too many and not focussed
- Data mining should be an *interactive* process
- The user directs what should be mined using a data mining query language or a graphical user interface
- Constraint-based mining
  - User flexibility: provides constraints on what to be mined (and what not)
  - System optimisation: explores such constraints for efficient mining

## Constraints in data mining

- Knowledge type constraint
  - Correlation, association, etc.
- Data constraint (use SQL like queries)
  - For example: *Find product pairs sold frequently in both stores in Sydney and Melbourne*
- Dimension / level constraint
  - In relevance to region, price, brand, customer category, etc.
- Rule or pattern constraint
  - Small sales (price < \$10) trigger big sales (sum > \$200)
- Interestingness constraint
  - Strong rules only:  $support_{min} > 3\%$ ,  $confidence_{min} > 75\%$

## Interestingness measure: Correlation (lift)

- Example: *Play basketball  $\Rightarrow$  Eat cereal [40%, 66.67%]* is misleading
  - If overall 75 % of all students eat cereal
  - *Play basketball  $\Rightarrow$  Not eat cereal [20%, 33.3%]* is more accurate, although with lower support and confidence
- Measure of dependent / correlated events: *Lift*

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

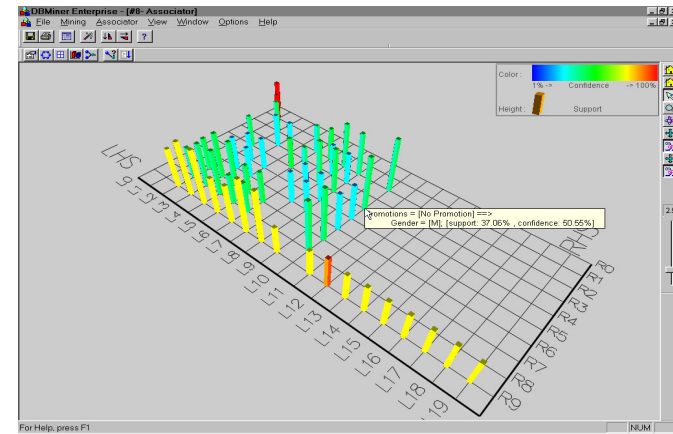
	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

$$lift(B, C) = \frac{2000/5000}{3000/5000 * 3750/5000} = 0.89 \quad lift(B, \neg C) = \frac{1000/5000}{3000/5000 * 1250/5000} = 1.33$$

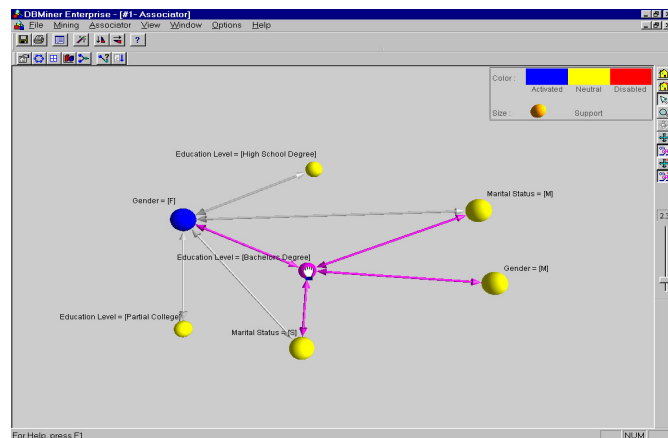
## More interestingness measures (Tan et al. 2002)

symbol	measure	range	formula
$\phi$	$\phi$ -coefficient	-1...1	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
$Q$	Yule's Q	-1...1	$\frac{P(A,B)P(A,B) - P(A)P(B)P(A,B)}{P(A,B)P(A,B) + P(A)P(B)P(A,B)}$
$Y$	Yule's Y	-1...1	$\frac{\sqrt{P(A,B)P(A,B)} - \sqrt{P(A)P(B)P(A,B)}}{\sqrt{P(A,B)P(A,B)} + \sqrt{P(A)P(B)P(A,B)}}$
$k$	Cohen's $k$	-1...1	$\frac{P(A,B) + P(A,B) - P(A)P(B) - P(A)P(B)}{1 - P(A)P(B) - P(A)P(B)}$
$PS$	Platietsky-Shapiro's $PS$	-0.25...0.25	$P(A,B) - P(A)P(B)$
$F$	Certainty factor	-1...1	$\max(\frac{P(A B) - P(A)}{1 - P(A)}, \frac{P(A B) - P(A)}{1 - P(A)})$
$AV$	add value	-0.5...1	$\max(P(B A) - P(B), P(A B) - P(A))$
$K$	Klorgen's $K$	-0.33...0.38	$\frac{P(A,B)}{P(A)P(B)} \max(P(A B) - P(B), P(A B) - P(A))$
$g$	Goodman-Kruskal's $g$	0...1	$\frac{\sum_{i,j} \max(P_{ij} - \min(P_{i.}, P_{.j}), 0)}{\sum_{i,j} \max(P_{ij} - \min(P_{i.}, P_{.j}), 0)}$
$M$	Mutual Information	0...1	$\sum_{i,j} P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}$
$J$	J-Measure	0...1	$\min(-\sum_i P(A_i) \log P(A_i), \log P(A)) - \sum_i P(A_i) \log P(A_i) \log P(B_i)$
			$\max(P(A) \log \frac{P(A)}{P(A)^2} + P(A) \log \frac{P(A)}{P(A)^2})$
$G$	Gini index	0...1	$\frac{P(A, B) \log \frac{P(A, B)}{P(A, B)} + P(A) \log \frac{P(A)}{P(A)^2}}{P(A, B) \log \frac{P(A, B)}{P(A, B)} + P(A) \log \frac{P(A)}{P(A)^2}}$
			$\max(P(A)(P(A)^2 + P(A)^2) + P(A)(P(A)^2 + P(A)(P(A)^2 - P(A)^2 - P(A)^2) - P(A)^2 - P(A)^2)$
$s$	support	0...1	$\frac{P(A, B)}{P(A)P(B)}$
$c$	confidence	0...1	$\max(P(B A), P(A B))$
$L$	Laplace	0...1	$\max(\frac{P(A B)+1}{N(P(A)+2)}, \frac{N(P(A)+1)}{N(P(A)+2)})$
$IS$	Cosine	0...1	$\frac{\sqrt{P(A)P(B)}}{\sqrt{P(A)P(B)}}$
$\gamma$	coherence(Jaccard)	0...1	$\frac{P(A) \cap P(B)}{P(A) \cup P(B)}$
$\alpha$	all confidence	0...1	$\max(\frac{P(A) \cap P(B)}{P(A) \cup P(B)}, \frac{P(A) \cap P(B)}{P(A) \cup P(B)})$
$o$	odds ratio	0... $\infty$	$\frac{P(A)P(B)}{P(A)P(B)}$
$V$	Conviction	0.5... $\infty$	$\max(\frac{P(A)P(B)}{P(A)P(B)}, \frac{P(B)P(A)}{P(A)P(B)})$
$\lambda$	lift	0... $\infty$	$\frac{P(A, B)}{P(A)P(B)}$
$S$	Collective strength	0... $\infty$	$\frac{P(A, B) + P(A, B)}{P(A)P(B) + P(A)P(B)} \times \frac{(1 - P(A)P(B) - P(A)P(B))}{1 - P(A)P(B) - P(A)P(B)}$
$\chi^2$	$\chi^2$	0... $\infty$	$\sum_i \frac{(P_{ij} - E_{ij})^2}{E_{ij}}$

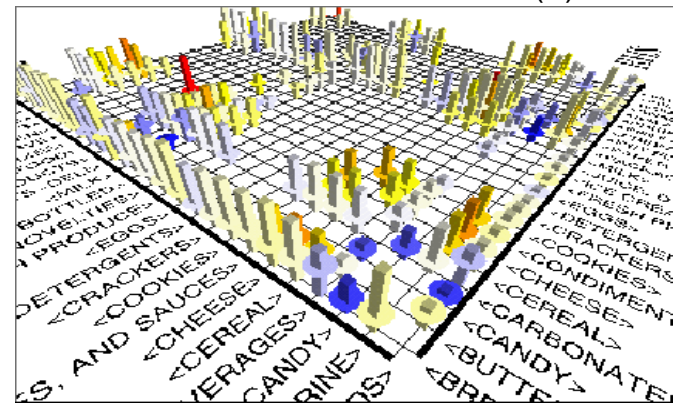
## Visualisation of association rules



## Visualisation of association rules (2)



### Visualisation of association rules (3)



## Review question

- Association rule mining allows us to predict what items customers will buy frequently together in a supermarket.

Yes    or    No?

- The rule: {'beer','chips'} → {'sausages'} [s=20%,c=40%] tells us that 40% of customers bought beer, chips and sausages.

Yes    or    No?

## What now.. things to do

- Read Chapter 6 and Sections 7.1-7.3 in text book
- Lab 2 next week: Read before you go into lab!  
Topic: Association mining - theoretical and in Rattle