



# Mining Association Rules

—From Association Mining to Correlation Analysis—

徐华

清华大学 计算机系 智能技术与系统国家重点实验室

xuhua@tsinghua.edu.cn

1

## Association and Correlations



- Association and Correlations
- Efficient and Scalable Frequent Itemset Mining Methods
- Mining Various Kinds of Association Rules
- From Association Mining to Correlation Analysis
- Constraint-based Association Mining

2



### Interestingness Measure: Correlations (Lift)



- ◉ *play basketball*  $\Rightarrow$  *eat cereal* (谷物) [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**

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

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

3



### Are All the Rules Found Interesting?



- ◉ “*Buy walnuts* (胡桃)  $\Rightarrow$  *buy milk* [1%, 80%]” is misleading
  - ◆ if 85% of customers buy milk
- ◉ Support and confidence are not good to represent correlations
- ◉ So many interestingness measures? (Tan, Kumar, Sritastava @KDD' 02)

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

$$all\_conf = \frac{\sup(X)}{\max\_item\_sup(X)}$$

$$coh = \frac{\sup(X)}{|universe(X)|}$$

	Milk	No Milk	Sum (row)
Coffee	m, c	$\sim m, c$	c
No Coffee	m, $\sim c$	$\sim m, \sim c$	$\sim c$
Sum(col.)	m	$\sim m$	$\Sigma$

DB	m, c	$\sim m, c$	m $\sim c$	$\sim m \sim c$	lift	all-conf	coh	$\chi^2$
A1	1000	100	100	10,000	9.26	0.91	0.83	9055
A2	100	1000	1000	100,000	8.44	0.09	0.05	670
A3	1000	100	10000	100,000	9.18	0.09	0.09	8172
A4	1000	1000	1000	1000	1	0.5	0.33	0

4



## Mining Highly Correlated Patterns



- *lift* and  $\chi^2$  are not good measures for correlations in transactional DBs
- *all-conf* or *coherence* could be good measures (Omiecinski @TKDE' 03)
- Both *all-conf* and *coherence* have the downward closure property
- Efficient algorithms can be derived for mining (Lee et al. @ICDM' 03sub)

$$all\_conf = \frac{\sup(X)}{\max\_item\_sup(X)}$$

$$coh = \frac{\sup(X)}{|universe(X)|}$$

DB	m, c	$\sim m, c$	$m \sim c$	$\sim m \sim c$	lift	all-conf	coh	$\chi^2$
A1	1000	100	100	10,000	9.26	0.91	0.83	9055
A2	100	1000	1000	100,000	8.44	0.09	0.05	670
A3	1000	100	10000	100,000	9.18	0.09	0.09	8172
A4	1000	1000	1000	1000	1	0.5	0.33	0

5



## Which Measures Should Be Used?



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	$\frac{P(A, B)P(\bar{A}, \bar{B}) - P(A, \bar{B})P(\bar{A}, B)}{P(A, B)P(\bar{A}, \bar{B}) + P(A, \bar{B})P(\bar{A}, B)}$
Y	Yule's Y	-1 ... 1	$\frac{P(A, B)P(\bar{A}, \bar{B}) - \sqrt{P(A, \bar{B})P(\bar{A}, B)}}$
k	Cohen's	-1 ... 1	$\frac{\sqrt{P(A, B)P(\bar{A}, \bar{B})} + \sqrt{P(A, \bar{B})P(\bar{A}, B)}}{P(A, B) + P(\bar{A}, \bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}$
PS	Piatetsky-Shapiro's	-0.25 ... 0.25	$\frac{P(A, B) - P(A)P(B)}{1 - P(A)P(B)}$
F	Certainty factor	-1 ... 1	$\max(\frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)})$
AV	added value	-0.5 ... 1	$\max(P(B A) - P(B), P(A B) - P(A))$
K	Klosgen's Q	0.33 ... 0.38	$\frac{\sqrt{P(\bar{A}, \bar{B})} \max(P(B A) - P(B), P(A B) - P(A))}{\sum_j \max_k P(A_j, B_k) + \sum_k \max_j P(A_j, B_k) - \max_k P(A_j) - \max_k P(B_k)}$
g	Goodman-kruskal's	0 ... 1	$\frac{2 - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
M	Mutual Information	0 ... 1	$\frac{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\sum_i \sum_j P(A_i) \log P(A_i) + \sum_i \sum_j P(B_j) \log P(B_j)}$
J	J-Measure	0 ... 1	$\frac{\min(-\sum_i P(A_i) \log P(A_i), -\sum_j P(B_j) \log P(B_j))}{\max(P(A, B) \log(\frac{P(A, B)}{P(A)P(B)}), P(A, B) \log(\frac{P(A, B)}{P(A)P(B)})}$
G	Gini index	0 ... 1	$\frac{P(A, B) \log \frac{P(A, B)}{P(A)P(B)} + P(\bar{A}, \bar{B}) \log \frac{P(\bar{A}, \bar{B})}{P(\bar{A})P(\bar{B})}}{P(A, B) \log \frac{P(A, B)}{P(A)P(B)} + P(\bar{A}, \bar{B}) \log \frac{P(\bar{A}, \bar{B})}{P(\bar{A})P(\bar{B})}}$
s	support	0 ... 1	$\max(P(A B), P(A B))$
c	confidence	0 ... 1	$\max(\frac{N P(A, B) + 1}{N P(A) + 2}, \frac{N P(A, B) + 1}{N P(B) + 2})$
L	Laplace	0 ... 1	$\frac{P(A, B)}{P(A)P(B)}$
IS	Cosine	0 ... 1	$\frac{P(A, B)}{\sqrt{P(A)P(B)}}$
$\gamma$	coherence (Jaccard)	0 ... 1	$\frac{P(A, B)}{P(A) + P(B) - P(A, B)}$
$\alpha$	all-confidence	0 ... 1	$\frac{\max(P(A B), P(B A))}{P(A, B)}$
o	odds ratio	0 ... $\infty$	$\frac{P(A, B)P(\bar{A}, \bar{B})}{P(A, \bar{B})P(\bar{A}, B)}$
V	Conviction	0.5 ... $\infty$	$\max(\frac{P(A, B)}{P(A)P(B)}, \frac{P(B, A)}{P(A)P(B)})$
$\lambda$	lift	0 ... $\infty$	$\frac{P(A, B)}{P(A)P(B)}$
S	Collective strength	0 ... $\infty$	$\frac{P(A, B) + P(\bar{A}, \bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A, B) - P(\bar{A}, \bar{B})}$
$\chi^2$	$\chi^2$	0 ... $\infty$	$\frac{P(A, B) - E_{AB}}{E_{AB}}$

6





**Thanks !**

