

Pattern Evaluation

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

Probability-based Measures

In a rule of the form $A \Rightarrow B$

- Support = $P(AB)$
- Confidence = $P(B|A)$
- Interest = $P(AB)/P(A)P(B)$
- Implication Strength = $P(A)P(\sim B)/P(A\sim B)$

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	\bar{y}	
X	f_{11}	f_{10}	f_{1+}
\bar{X}	f_{01}	f_{00}	f_{0+}
	f_{+1}	f_{+0}	$ T $

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

Used to define various measures

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

Drawback of Confidence

	Coffee	$\bar{\text{Coffee}}$	
Tea	15	5	20
$\bar{\text{Tea}}$	75	5	80
	90	10	100

Association Rule: Tea \rightarrow Coffee

Confidence = $P(\text{Coffee}|\text{Tea}) = 0.75$

but $P(\text{Coffee}) = 0.9$

\Rightarrow Although confidence is high, rule is misleading

$\Rightarrow P(\text{Coffee}|\bar{\text{Tea}}) = 0.9375$

Criticism to Support and Confidence

Example 1: (Aggarwal & Yu, paper at PODS98)

- Among 1000 students
 - 600 play basketball
 - 750 eat cereal
 - 400 both play basketball and eat cereal
- *play basketball \Rightarrow eat cereal* [40%, 66.7%] is misleading because 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 far more accurate, although with lower support and confidence

	basketball	not basketball	sum(row)
cereal	400	350	750
not cereal	200	50	250
sum(col.)	600	400	1000

Criticism to Support and Confidence (Cont.)

Example 2:

- X and Y : positively correlated
- X and Z : negatively related
- support and confidence of $X \Rightarrow Z$ dominates

X	1	1	1	0	0	0	0
Y	1	1	0	0	0	0	0
Z	0	1	1	1	1	1	1

Measure of dependence or correlation of events

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

Rule	Support	Confidence
$X \Rightarrow Y$	25%	50%
$X \Rightarrow Z$	37.50%	75%

$P(B|A)/P(B)$ is also called the **lift** of rule $A \Rightarrow B$

Other Interestingness Measures: Interest

- Interest (correlation, lift) $\frac{P(A \wedge B)}{P(A)P(B)}$
 - taking both P(A) and P(B) in consideration
 - $P(A \wedge B) = P(B) \cdot P(A)$, if A and B are independent events
 - A and B negatively correlated, if the value is less than 1; otherwise A and B positively correlated

X	1	1	1	1	0	0	0	0
Y	1	1	0	0	0	0	0	0
Z	0	1	1	1	1	1	1	1

Itemset	Support	Interest
X,Y	25%	2
X,Z	37.50%	0.9
Y,Z	12.50%	0.57

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)
- $P(S \wedge B) = 420/1000 = 0.42$
- $P(S) \times P(B) = 0.6 \times 0.7 = 0.42$
- $P(S \wedge B) = P(S) \times P(B) \Rightarrow$ Statistical independence
- $P(S \wedge B) > P(S) \times P(B) \Rightarrow$ Positively correlated
- $P(S \wedge B) < P(S) \times P(B) \Rightarrow$ Negatively correlated

Statistical-based Measures

- Measures that take into account statistical dependence

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

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

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

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

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 = $P(\text{Coffee} | \text{Tea}) = 0.75$

but $P(\text{Coffee}) = 0.9$

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

Drawback of Lift & Interest

	y	\overline{y}	
x	10	0	10
\overline{x}	0	90	90
	10	90	100

$$Lift = \frac{0.1}{(0.1)(0.1)} = 10$$

	y	\overline{y}	
x	90	0	90
\overline{x}	0	10	10
	90	10	100

$$Lift = \frac{0.9}{(0.9)(0.9)} = 1.11$$

Statistical independence:

If $P(X, Y) = P(X)P(Y) \Rightarrow$ Lift = 1

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
- $M(A, B)$ increases 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

Property under Variable Permutation

	B	\bar{B}
A	p	q
\bar{A}	r	s

 \implies

	B	\bar{B}
A	p	q
\bar{A}	r	s

Does $M(A,B) = M(B,A)$?

Symmetric measures:

- ◆ support, lift, collective strength, cosine, Jaccard, etc

Asymmetric measures:

- ◆ confidence, conviction, Laplace, J-measure, etc

Property under Row/Column Scaling

Grade-Gender Example (Mosteller, 1968):

	Male	Female	
High	2	3	5
Low	1	4	5
	3	7	10

	Male	Female	
High	4	30	34
Low	2	40	42
	6	70	76

\downarrow \downarrow
 2x 10x

Mosteller:

Underlying association should be independent of the relative number of male and female students in the samples

Property under Inversion Operation

	A	B		C	D		E	F
Transaction 1	1	0		0	1		0	0
■	0	0		1	1		1	0
■	0	0		1	1		1	0
■	0	1		1	0		1	1
■	0	0		1	1		1	0
■	0	0		1	1		1	0
■	0	0		1	1		1	0
■	0	0		1	1		1	0
Transaction N	1	0		0	1		0	0

(a) (b) (c)

Example: ϕ -Coefficient

• ϕ -coefficient is analogous to correlation coefficient for continuous variables

	y	\bar{y}	
X	60	10	70
\bar{X}	10	20	30
	70	30	100

	y	\bar{y}	
X	20	10	30
\bar{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

Property under Null Addition

	B	\bar{B}
A	p	q
\bar{A}	r	s

 \implies

	B	\bar{B}
A	p	q
\bar{A}	r	s+k

Invariant measures:

- ◆ support, cosine, Jaccard, etc

Non-invariant measures:

- ◆ correlation, Gini, mutual information, odds ratio, etc

Interestingness Measurements

Objective measures

Two popular measurements:

- *support*; and
- *confidence*

Subjective measures [Silberschatz & Tuzhilin, KDD95]

A rule (pattern) is interesting if

- it is *unexpected* (surprising to the user); and/or
- *actionable* (the user can do something with it)