# Pattern Evaluation
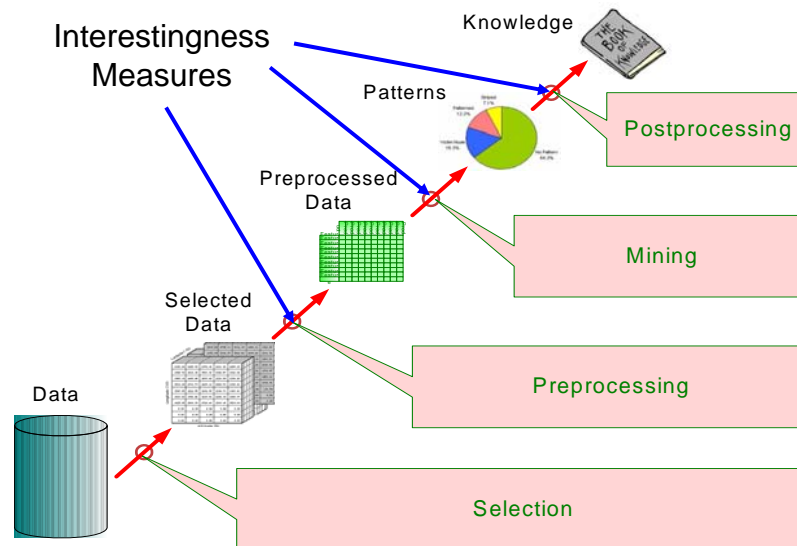
- Association rule algorithms tend to produce too many rules
  - many of them are uninteresting or redundant
  - Redundant if {A,B,C} $\rightarrow$ {D} and {A,B} $\rightarrow$ {D} have same support & confidence

- 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

# Application of Interestingness Measure



Interestingness Measures

Knowledge

Patterns

Preprocessed Data

Selected Data

Data

Postprocessing

Mining

Preprocessing

Selection

# 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 $\underline{X}$ and $\overline{Y}$
$f_{01}$: support of $\overline{X}$ and $\underline{Y}$
$f_{00}$: support of $\overline{X}$ and $\overline{Y}$

Used to define various measures

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

---

# Drawback of Confidence

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

Association Rule: Tea $\rightarrow$ Coffee

Confidence= P(Coffee|Tea) = 0.75

but P(Coffee) = 0.9

$\Rightarrow$ Although confidence is high, rule is misleading

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

# 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

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

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

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

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

# Example: Lift/Interest

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

Association Rule: Tea → Coffee

Confidence= P(Coffee|Tea) = 0.75

but P(Coffee) = 0.9

⇒ 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 |

|      | Y  | $\overline{Y}$ |     |
|------|----|----|-----|
| X    | 90 | 0  | 90  |
| $\overline{X}$ | 0  | 10 | 10  |
|      | 90 | 10 | 100 |

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

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

**Statistical independence:**

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

**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 | $\frac{P(A,B)-P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$ |
| 2 | Goodman-Kruskal's $(\lambda)$ | $\frac{\sum_j \max_k P(A_j,B_k)+\sum_k \max_j P(A_j,B_k)-\max_j P(A_j)-\max_k P(B_k)}{2-\max_j P(A_j)-\max_k P(B_k)}$ |
| 3 | Odds ratio $(\alpha)$ | $\frac{P(A,B)P(\overline{A},\overline{B})}{P(A,\overline{B})P(\overline{A},B)}$ |
| 4 | Yule's $Q$ | $\frac{P(A,B)P(\overline{A}\overline{B})-P(A,\overline{B})P(\overline{A},B)}{P(A,B)P(\overline{A}\overline{B})+P(A,\overline{B})P(\overline{A},B)}=\frac{\alpha-1}{\alpha+1}$ |
| 5 | Yule's $Y$ | $\frac{\sqrt{P(A,B)P(\overline{A}\overline{B})}-\sqrt{P(A,\overline{B})P(\overline{A},B)}}{\sqrt{P(A,B)P(\overline{A}\overline{B})}+\sqrt{P(A,\overline{B})P(\overline{A},B)}}=\frac{\sqrt{\alpha}-1}{\sqrt{\alpha}+1}$ |
| 6 | Kappa $(\kappa)$ | $\frac{P(A,B)+P(\overline{A},\overline{B})-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A)P(B)-P(\overline{A})P(\overline{B})}$ |
| 7 | Mutual Information $(M)$ | $\frac{\sum_i \sum_j P(A_i,B_j)\log\frac{P(A_i,B_j)}{P(A_i)P(B_j)}}{\min(-\sum_i P(A_i)\log P(A_i),-\sum_j P(B_j)\log P(B_j))}$ |
| 8 | J-Measure $(J)$ | $\max\left(P(A,B)\log(\frac{P(B|A)}{P(B)})+P(A\overline{B})\log(\frac{P(\overline{B}|A)}{P(\overline{B})}),\right.$ $\left. P(A,B)\log(\frac{P(A|B)}{P(A)})+P(\overline{A}B)\log(\frac{P(\overline{A}|B)}{P(A)})\right)$ |
| 9 | Gini index $(G)$ | $\max\left(P(A)[P(B|A)^2+P(\overline{B}|A)^2]+P(\overline{A})[P(B|\overline{A})^2+P(\overline{B}|\overline{A})^2]\right.$ $-P(B)^2-P(\overline{B})^2,$ $P(B)[P(A|B)^2+P(\overline{A}|B)^2]+P(\overline{B})[P(A|\overline{B})^2+P(\overline{A}|\overline{B})^2]$ $\left.-P(A)^2-P(\overline{A})^2\right)$ |
| 10 | Support $(s)$ | $P(A,B)$ |
| 11 | Confidence $(c)$ | $\max(P(B|A),P(A|B))$ |
| 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 $(I)$ | $\frac{P(A,B)}{P(A)P(B)}$ |
| 15 | cosine $(IS)$ | $\frac{P(A,B)}{\sqrt{P(A)P(B)}}$ |
| 16 | Piatetsky-Shapiro's $(PS)$ | $P(A,B)-P(A)P(B)$ |
| 17 | Certainty factor $(F)$ | $\max\left(\frac{P(B|A)-P(B)}{1-P(B)},\frac{P(A|B)-P(A)}{1-P(A)}\right)$ |
| 18 | Added Value $(AV)$ | $\max(P(B|A)-P(B),P(A|B)-P(A))$ |
| 19 | Collective strength $(S)$ | $\frac{P(A,B)+P(\overline{AB})}{P(A)P(B)+P(\overline{A})P(\overline{B})}\times\frac{1-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A,B)-P(\overline{AB})}$ |
| 20 | Jaccard $(\zeta)$ | $\frac{P(A,B)}{P(A)+P(B)-P(A,B)}$ |
| 21 | Klosgen $(K)$ | $\sqrt{P(A,B)}\max(P(B|A)-P(B),P(A|B)-P(A))$ |

---

# Comparing Different Measures

10 examples of contingency tables:

| Example | $f_{11}$ | $f_{10}$ | $f_{01}$ | $f_{00}$ |
|---------|------|------|------|------|
| E1 | 8123 | 83 | 424 | 1370 |
| E2 | 8330 | 2 | 622 | 1046 |
| E3 | 9481 | 94 | 127 | 298 |
| E4 | 3954 | 3080 | 5 | 2961 |
| E5 | 2886 | 1363 | 1320 | 4431 |
| E6 | 1500 | 2000 | 500 | 6000 |
| E7 | 4000 | 2000 | 1000 | 3000 |
| E8 | 4000 | 2000 | 2000 | 2000 |
| E9 | 1720 | 7121 | 5 | 1154 |
| E10 | 61 | 2483 | 4 | 7452 |

Rankings of contingency tables using various measures:

| # | $\phi$ | $\lambda$ | $\alpha$ | $Q$ | $Y$ | $\kappa$ | $M$ | $J$ | $G$ | $s$ | $c$ | $L$ | $V$ | $I$ | $IS$ | $PS$ | $F$ | $AV$ | $S$ | $\zeta$ | $K$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| E1 | 1 | 1 | 3 | 3 | 3 | 1 | 2 | 2 | 1 | 3 | 5 | 5 | 4 | 6 | 2 | 2 | 4 | 6 | 1 | 2 | 5 |
| E2 | 2 | 2 | 1 | 1 | 1 | 2 | 1 | 3 | 2 | 2 | 1 | 1 | 1 | 8 | 3 | 5 | 1 | 8 | 2 | 3 | 6 |
| E3 | 3 | 3 | 4 | 4 | 4 | 3 | 3 | 8 | 7 | 1 | 4 | 4 | 6 | 10 | 1 | 8 | 6 | 10 | 3 | 1 | 10 |
| E4 | 4 | 7 | 2 | 2 | 2 | 5 | 4 | 1 | 3 | 6 | 2 | 2 | 2 | 4 | 4 | 1 | 2 | 3 | 4 | 5 | 1 |
| E5 | 5 | 4 | 8 | 8 | 8 | 4 | 7 | 5 | 4 | 7 | 9 | 9 | 9 | 3 | 6 | 3 | 9 | 4 | 5 | 6 | 3 |
| E6 | 6 | 6 | 7 | 7 | 7 | 7 | 6 | 4 | 6 | 9 | 8 | 8 | 7 | 2 | 8 | 6 | 7 | 2 | 7 | 8 | 2 |
| E7 | 7 | 5 | 9 | 9 | 9 | 6 | 8 | 6 | 5 | 4 | 7 | 7 | 8 | 5 | 5 | 4 | 8 | 5 | 6 | 4 | 4 |
| E8 | 8 | 9 | 10 | 10 | 10 | 8 | 10 | 10 | 8 | 4 | 10 | 10 | 10 | 9 | 7 | 7 | 10 | 9 | 8 | 7 | 9 |
| E9 | 9 | 9 | 5 | 5 | 5 | 9 | 9 | 7 | 9 | 8 | 3 | 3 | 3 | 7 | 9 | 9 | 3 | 7 | 9 | 9 | 8 |
| E10 | 10 | 8 | 6 | 6 | 6 | 10 | 5 | 9 | 10 | 10 | 6 | 6 | 5 | 1 | 10 | 10 | 5 | 1 | 10 | 10 | 7 |

# 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** | **B̄** |
|-----|-------|-------|
| **A** | p | q |
| **Ā** | r | s |

$\Rightarrow$

|     | **A** | **Ā** |
|-----|-------|-------|
| **B** | p | r |
| **B̄** | q | 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 |

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 | 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 | $\overline{Y}$ |   |
|---|---|---|---|
| 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**

---

# Property under Null Addition

|   | **B** | **$\overline{B}$** |
|---|---|---|
| **A** | p | q |
| **$\overline{A}$** | r | s |

$\Longrightarrow$

|   | **B** | **$\overline{B}$** |
|---|---|---|
| **A** | p | q |
| **$\overline{A}$** | r | s + k |

Invariant measures:

- ◆ support, cosine, Jaccard, etc

Non-invariant measures:

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

## Different Measures have Different Properties

| Symbol | Measure | Range | P1 | P2 | P3 | O1 | O2 | O3 | O3' | O4 |
|---|---|---|---|---|---|---|---|---|---|---|
| $\Phi$ | Correlation | -1 ... 0 ... 1 | Yes | Yes | Yes | Yes | No | Yes | Yes | No |
| $\lambda$ | Lambda | 0 ... 1 | Yes | No | No | Yes | No | No* | Yes | No |
| $\alpha$ | Odds ratio | 0 ... 1 ... $\infty$ | Yes* | Yes | Yes | Yes | Yes | Yes* | Yes | No |
| Q | Yule's Q | -1 ... 0 ... 1 | Yes | Yes | Yes | Yes | Yes | Yes | Yes | No |
| Y | Yule's Y | -1 ... 0 ... 1 | Yes | Yes | Yes | Yes | Yes | Yes | Yes | No |
| $\kappa$ | Cohen's | -1 ... 0 ... 1 | Yes | Yes | Yes | Yes | No | No | Yes | No |
| M | Mutual Information | 0 ... 1 | Yes | Yes | Yes | Yes | No | No* | Yes | No |
| J | J-Measure | 0 ... 1 | Yes | No | No | No | No | No | No | No |
| G | Gini Index | 0 ... 1 | Yes | No | No | No | No | No* | Yes | No |
| s | Support | 0 ... 1 | No | Yes | No | Yes | No | No | No | No |
| c | Confidence | 0 ... 1 | No | Yes | No | Yes | No | No | No | Yes |
| L | Laplace | 0 ... 1 | No | Yes | No | Yes | No | No | No | No |
| V | Conviction | 0.5 ... 1 ... $\infty$ | No | Yes | No | Yes** | No | No | Yes | No |
| I | Interest | 0 ... 1 ... $\infty$ | Yes* | Yes | Yes | Yes | No | No | No | No |
| IS | IS (cosine) | 0 .. 1 | No | Yes | Yes | Yes | No | No | No | Yes |
| PS | Piatetsky-Shapiro's | -0.25 ... 0 ... 0.25 | Yes | Yes | Yes | Yes | No | Yes | Yes | No |
| F | Certainty factor | -1 ... 0 ... 1 | Yes | Yes | Yes | No | No | No | Yes | No |
| AV | Added value | 0.5 ... 1 ... 1 | Yes | Yes | Yes | No | No | No | No | No |
| S | Collective strength | 0 ... 1 ... $\infty$ | No | Yes | Yes | Yes | No | Yes* | Yes | No |
| $\zeta$ | Jaccard | 0 .. 1 | No | Yes | Yes | Yes | No | No | No | Yes |
| K | Klosgen's | $\left(\sqrt{\frac{2}{\sqrt{3}}}-1\right)\left(2-\sqrt{3}-\frac{1}{\sqrt{3}}\right)...0...\frac{2}{3\sqrt{3}}$ | Yes | Yes | Yes | No | No | No | No | No |

---

## Support-based Pruning

- Most of the association rule mining algorithms use support measure to prune rules and itemsets

- Study effect of support pruning on correlation of itemsets
  - Generate 10000 random contingency tables
  - Compute support and pairwise correlation for each table
  - Apply support-based pruning and examine the tables that are removed
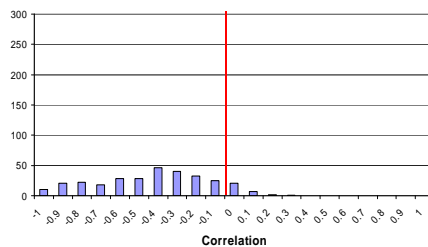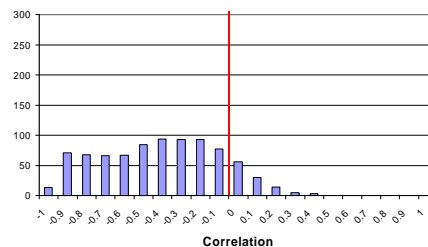
# Effect of Support-based Pruning

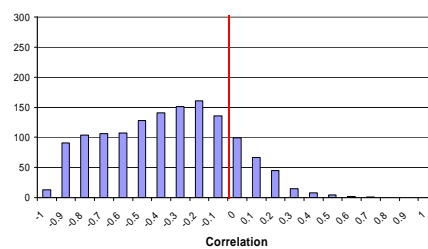**All Itempairs**


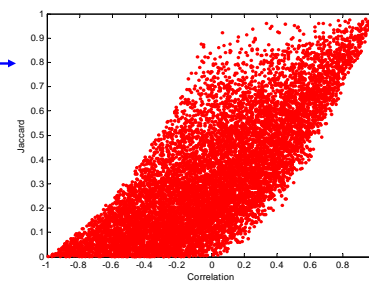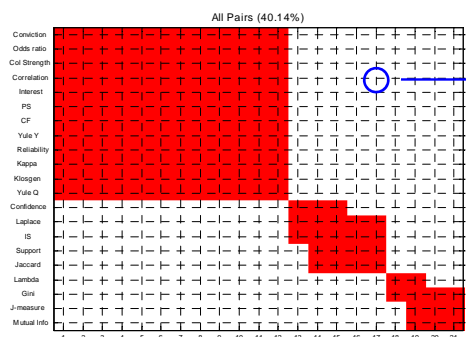
Correlation

# Effect of Support-based Pruning



Support-based pruning eliminates mostly negatively correlated itemsets

# Effect of Support-based Pruning

- Investigate how support-based pruning affects other measures

- Steps:
  – Generate 10000 contingency tables
  – Rank each table according to the different measures
  – Compute the pair-wise correlation between the measures

---

# Effect of Support-based Pruning

- Without Support Pruning (All Pairs)



Scatter Plot between Correlation & Jaccard Measure

- Red cells indicate correlation between the pair of measures > 0.85

- 40.14% pairs have correlation > 0.85

# Effect of Support-based Pruning

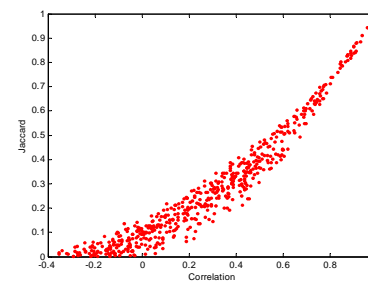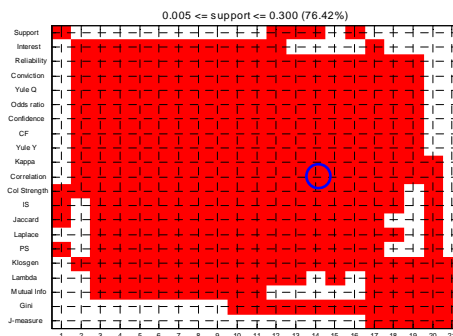- ◆ 0.5% ≤ support ≤ 50%



0.005 <= support <= 0.500 (61.45%)

Scatter Plot between Correlation & Jaccard Measure:

- ◆ 61.45% pairs have correlation > 0.85

---

# Effect of Support-based Pruning

- ◆ 0.5% ≤ support ≤ 30%



0.005 <= support <= 0.300 (76.42%)

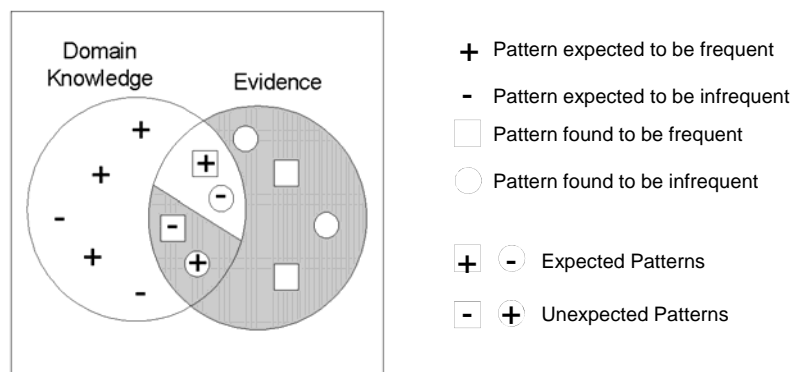Scatter Plot between Correlation & Jaccard Measure

- ◆ 76.42% pairs have correlation > 0.85

# Subjective Interestingness Measure

- Objective measure:
  - Rank patterns based on statistics computed from data
  - e.g., 21 measures of association (support, confidence, Laplace, Gini, mutual information, Jaccard, etc).

- Subjective measure:
  - Rank patterns according to user's interpretation
    - A pattern is subjectively interesting if it contradicts the expectation of a user (Silberschatz & Tuzhilin)
    - A pattern is subjectively interesting if it is actionable (Silberschatz & Tuzhilin)

---

# Interestingness via Unexpectedness

- Need to model expectation of users (domain knowledge)



+ Pattern expected to be frequent

- Pattern expected to be infrequent

☐ Pattern found to be frequent

◯ Pattern found to be infrequent

+ ◯ Expected Patterns

- + Unexpected Patterns

- Need to combine expectation of users with evidence from data (i.e., extracted patterns)

# Simpson's Paradox

- Hidden influences can produce misleading results
  - C(EM=Y | HDTV=Y) = 55%
  - C(EM=Y | HDTV=N) = 45%
  - Conclude HDTV -> EM?

| Buy HDTV | Buy Exercise Machine | | |
|---|---|---|---|
| | Yes | No | |
| Yes | 99 | 81 | 180 |
| No | 54 | 66 | 120 |
| Totals | 153 | 147 | 300 |

---

# Simpson's Paradox

College Students

C(EM=Y | HDTV=Y) = 10%

C(EM=Y | HDTV=N) = 11.8%

Conclude HDTV -> not EM

Working Adults

C(EM=Y | HDTV=Y) = 57.7%

C(EM=Y | HDTV=N) = 58.1%

Conclude HDTV -> not EM

| Customer Group | Buy HDTV | Buy Exercise Machine | | Total |
|---|---|---|---|---|
| | | Yes | No | |
| College Students | Yes | 1 | 9 | 10 |
| | No | 4 | 30 | 34 |
| Working Adults | Yes | 98 | 72 | 170 |
| | No | 50 | 36 | 86 |