

The background of the slide is a complex, abstract composition. It features a network of thin, reddish-brown lines connecting various points, creating a web-like structure. Interspersed among these lines are numerous small, colored dots in shades of green, blue, and orange. The overall aesthetic is technical and data-driven, with a mix of organic and geometric forms. A large, white, angular shape serves as a backdrop for the title text.

Comparison of Null-Invariant Measures

Comparison of Null-Invariant Measures

- ❑ Not all null-invariant measures are created equal
- ❑ Which one is better?
 - ❑ $D_4 - D_6$ differentiate the null-invariant measures
 - ❑ Kulc (Kulczynski 1927) holds firm and is in balance of both directional implications

2-variable contingency table

| | <i>milk</i> | $\neg milk$ | Σ_{row} |
|----------------|-------------|----------------|----------------|
| <i>coffee</i> | <i>mc</i> | $\neg mc$ | <i>c</i> |
| $\neg coffee$ | $m\neg c$ | $\neg m\neg c$ | $\neg c$ |
| Σ_{col} | <i>m</i> | $\neg m$ | Σ |

All 5 are null-invariant

| Data set | <i>mc</i> | $\neg mc$ | $m\neg c$ | $\neg m\neg c$ | <i>AllConf</i> | Jaccard | <i>Cosine</i> | <i>Kulc</i> | <i>MaxConf</i> |
|----------|-----------|-----------|-----------|----------------|----------------|---------|---------------|-------------|----------------|
| D_1 | 10,000 | 1,000 | 1,000 | 100,000 | 0.91 | 0.83 | 0.91 | 0.91 | 0.91 |
| D_2 | 10,000 | 1,000 | 1,000 | 100 | 0.91 | 0.83 | 0.91 | 0.91 | 0.91 |
| D_3 | 100 | 1,000 | 1,000 | 100,000 | 0.09 | 0.05 | 0.09 | 0.09 | 0.09 |
| D_4 | 1,000 | 1,000 | 1,000 | 100,000 | 0.5 | 0.33 | 0.5 | 0.5 | 0.5 |
| D_5 | 1,000 | 100 | 10,000 | 100,000 | 0.09 | 0.09 | 0.29 | 0.5 | 0.91 |
| D_6 | 1,000 | 10 | 100,000 | 100,000 | 0.01 | 0.01 | 0.10 | 0.5 | 0.99 |

Subtle: They disagree on those cases

Imbalance Ratio with Kulczynski Measure

- IR (Imbalance Ratio): measure the imbalance of two itemsets A and B in rule implications:

$$IR(A, B) = \frac{|s(A) - s(B)|}{s(A) + s(B) - s(A \cup B)}$$

- Kulczynski and Imbalance Ratio (IR) together present a clear picture for all the three datasets D_4 through D_6
 - D_4 is neutral & balanced; D_5 is neutral but imbalanced
 - D_6 is neutral but very imbalanced

| Data set | mc | $\neg mc$ | $m\neg c$ | $\neg m\neg c$ | Jaccard | <i>Cosine</i> | <i>Kulc</i> | IR |
|----------|--------|-----------|-----------|----------------|---------|---------------|-------------|------|
| D_1 | 10,000 | 1,000 | 1,000 | 100,000 | 0.83 | 0.91 | 0.91 | 0 |
| D_2 | 10,000 | 1,000 | 1,000 | 100 | 0.83 | 0.91 | 0.91 | 0 |
| D_3 | 100 | 1,000 | 1,000 | 100,000 | 0.05 | 0.09 | 0.09 | 0 |
| D_4 | 1,000 | 1,000 | 1,000 | 100,000 | 0.33 | 0.5 | 0.5 | 0 |
| D_5 | 1,000 | 100 | 10,000 | 100,000 | 0.09 | 0.29 | 0.5 | 0.89 |
| D_6 | 1,000 | 10 | 100,000 | 100,000 | 0.01 | 0.10 | 0.5 | 0.99 |

Example: Analysis of DBLP Coauthor Relationships

Recent DB conferences, removing balanced associations, low sup, etc.

| ID | Author <i>A</i> | Author <i>B</i> | $s(A \cup B)$ | $s(A)$ | $s(B)$ | Jaccard | <i>Cosine</i> | <i>Kulc</i> |
|----|----------------------|----------------------|---------------|--------|--------|------------|---------------|-------------|
| 1 | Hans-Peter Kriegel | Martin Ester | 28 | 146 | 54 | 0.163 (2) | 0.315 (7) | 0.355 (9) |
| 2 | Michael Carey | Miron Livny | 26 | 104 | 58 | 0.191 (1) | 0.335 (4) | 0.349 (10) |
| 3 | Hans-Peter Kriegel | Joerg Sander | 24 | 146 | 36 | 0.152 (3) | 0.331 (5) | 0.416 (8) |
| 4 | Christos Faloutsos | Spiros Papadimitriou | 20 | 162 | 26 | 0.119 (7) | 0.308 (10) | 0.446 (7) |
| 5 | Hans-Peter Kriegel | Martin Pfeifle | 18 | 146 | 18 | 0.123 (6) | 0.351 (2) | 0.562 (2) |
| 6 | Hector Garcia-Molina | Wilburt Labio | 16 | 144 | 18 | 0.110 (9) | 0.314 (8) | 0.500 (4) |
| 7 | Divyakant Agrawal | Wang Hsiung | 16 | 120 | 16 | 0.133 (5) | 0.365 (1) | 0.567 (1) |
| 8 | Elke Rundensteiner | Murali Mani | 16 | 104 | 20 | 0.148 (4) | 0.351 (3) | 0.477 (6) |
| 9 | Divyakant Agrawal | Oliver Po | 12 | 120 | 12 | 0.100 (10) | 0.316 (6) | 0.550 (3) |
| 10 | Gerhard Weikum | Martin Theobald | 12 | 106 | 14 | 0.111 (8) | 0.312 (9) | 0.485 (5) |

Advisor-advisee relation: Kulc: high, Jaccard: low, cosine: middle

- ❑ Which pairs of authors are strongly related?
 - ❑ Use Kulc to find Advisor-advisee, close collaborators