# Proximity Measure for Binary Attributes

❑ A contingency table for binary data

<table>
<tr><td></td><td colspan="3" align="center">Object <em>j</em></td></tr>
<tr><td></td><td>1</td><td>0</td><td>sum</td></tr>
<tr><td rowspan="2">Object <em>i</em></td></tr>
<tr><td>1</td><td>$q$</td><td>$r$</td><td>$q+r$</td></tr>
<tr><td>0</td><td>$s$</td><td>$t$</td><td>$s+t$</td></tr>
<tr><td>sum</td><td>$q+s$</td><td>$r+t$</td><td>$p$</td></tr>
</table>

❑ Distance measure for symmetric binary variables  $d(i, j) = \dfrac{r+s}{q+r+s+t}$

❑ Distance measure for asymmetric binary variables:  $d(i, j) = \dfrac{r+s}{q+r+s}$

❑ Jaccard coefficient (*similarity* measure for *asymmetric* binary variables):  $sim_{Jaccard}(i, j) = \dfrac{q}{q+r+s}$

❑ Note: Jaccard coefficient is the same as  (a concept discussed in Pattern Discovery)

$$coherence(i, j) = \frac{sup(i, j)}{sup(i) + sup(j) - sup(i, j)} = \frac{q}{(q+r) + (q+s) - q}$$

62

# Example: Dissimilarity between Asymmetric Binary Variables

| Name | Gender | Fever | Cough | Test-1 | Test-2 | Test-3 | Test-4 |
|------|--------|-------|-------|--------|--------|--------|--------|
| Jack | M | Y | N | P | N | N | N |
| Mary | F | Y | N | P | N | P | N |
| Jim | M | Y | P | N | N | N | N |

- ❑ Gender is a symmetric attribute (not counted in)

- ❑ The remaining attributes are asymmetric binary

- ❑ Let the values Y and P be 1, and the value N be 0

- ❑ Distance: $d(i, j) = \dfrac{r + s}{q + r + s}$

$$d(jack, mary) = \frac{0+1}{2+0+1} = 0.33$$

$$d(jack, jim) = \frac{1+1}{1+1+1} = 0.67$$

$$d(jim, mary) = \frac{1+2}{1+1+2} = 0.75$$

|  | Mary 1 | 0 | Σ_row |
|------|------|------|------|
| Jack 1 | 2 | ~~0~~ 1 | 2 |
| 0 | 1 | 3 | 4 |
| Σ_col | 3 | ~~3~~ 4 | 6 |

|  | Jim 1 | 0 | Σ_row |
|------|------|------|------|
| Jack 1 | ~~1~~ 2 | 1 | ~~2~~ 3 |
| 0 | 1 | 3 | 4 |
| Σ_col | ~~2~~ 3 | 4 | ~~6~~ 7 |

|  | Mary 1 | 0 | Σ_row |
|------|------|------|------|
| Jim 1 | 1 | 1 | 2 |
| 0 | 2 | 2 | 4 |
| Σ_col | 3 | 3 | 6 |

# Proximity Measure for Categorical Attributes

❑ Categorical data, also called nominal attributes → เป็นหมวดหมู่

   ❑ Example: Color (red, yellow, blue, green), profession, etc.

❑ Method 1: Simple matching

   ❑ *m*: # of matches, *p*: total # of variables

$$d(i, j) = \frac{p - m}{p} \longrightarrow \frac{\text{ทั้งหมด} - \text{เหมือน}}{\text{ทั้งหมด}}$$

❑ Method 2: Use a large number of binary attributes

   ❑ Creating a new binary attribute for each of the *M* nominal states

# Ordinal Variables

❑ An ordinal variable can be discrete or continuous → เรียงลำดับ

❑ Order is important, e.g., rank (e.g., freshman, sophomore, junior, senior)

❑ Can be treated like interval-scaled

กลำดับที่ ?

  ❑ Replace *an ordinal variable value* by its rank:  $r_{if} \in \{1, ..., M_f\}$

  ❑ Map the range of each variable onto [0, 1] by replacing *i*-th object in the *f*-th variable by

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$  → freshman = $\frac{1-1}{4-1} = \frac{0}{3} = 0$

    ❑ Example:  freshman: 0; sophomore: 1/3; junior: 2/3; senior 1

      ❑ Then distance:  d(freshman, senior) = 1, d(junior, senior) = 1/3

  ❑ Compute the dissimilarity using methods for interval-scaled variables  $|1-0| = |\frac{2}{3} - \frac{3}{3}| = \frac{1}{3}$

65