# Proximity Measure for Binary Attributes

- A contingency table for binary data

|  | Object $j$ | | |
|---|---|---|---|
| Object $i$ | **1** | **0** | **sum** |
| 1 | $q$ | $r$ | $q+r$ |
| 0 | $s$ | $t$ | $s+t$ |
| sum | $q+s$ | $r+t$ | $p$ |

*ตอบถูก ขาย = ขาย*

*0 ตรงกับ 0
ไม่ใช่ = ไม่ใช่*

- Distance measure for <u>symmetric</u> binary variables $\quad d(i,j) = \dfrac{r+s}{q+r+s+\boxed{t}}$

- Distance measure for <u>asymmetric</u> binary variables: $\quad d(i,j) = \dfrac{r+s}{q+r+s}$

  *ขาย = แคใว่ว=?*

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

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

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

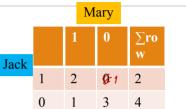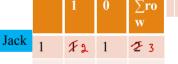# Example: Dissimilarity between Asymmetric Binary Variables

| Name | Gender | Fever | Cough | Test-1 | Test-2 | Test-3 | Test-4 |
|------|--------|-------|-------|--------|--------|--------|--------|
| Jack | M 1 | Y 1 | N 0 | P 1 | N 0 | N 0 | N 0 |
| Mary | F 0 | Y 1 | N 0 | P 1 | N 0 | P 1 | N 0 |
| Jim  | M 1 | Y 1 | P 1 | N 0 | N 0 | N 0 | N 0 |

- 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 / Jack table:

| | 1 | 0 | ∑row |
|---|---|---|---|
| 1 | 2 | 1 | 2 |
| 0 | 1 | 3 | 4 |
| ∑c | 3 | 4 | 7 |

Jim / Jack table:

| | 1 | 0 | ∑row |
|---|---|---|---|
| 1 | 1 | 2 | 1 | 3 |
| 0 | 1 | 3 | 4 |
| ∑c | 3 | 4 | 7 |

Mary / Jim table:

| | 1 | 0 | ∑row |
|---|---|---|---|
| 1 | 1 | 1 | 2 |
| 0 | 2 | 2 | 4 |

| | 1 | 0 | SUM |
|---|---|---|---|
| 1 | q | r | |
| 0 | s | t | |
| SUM | | | |

# Proximity Measure for Categorical Attributes

- Categorical data, also called nominal attributes *ที่มีหลายตัวอย่าง*

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

- <u>Method 1</u>: Simple matching

  - $m$: # of matches, $p$: total # of variables *จำนวนตัวที่ไม่เหมือน*

$$d(i, j) = \frac{p - m}{p}$$

*ตัวหมด - เหมือน*
*ตัวหมด*

- <u>Method 2</u>: 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)   *(1,2,3,4)*

- Can be treated like interval-scaled

  *ลำดับที่เท่าไหร่*

  - 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   $z_{if} = \dfrac{r_{if} - 1}{M_f - 1}$   *fresh man* $= \dfrac{1-1}{4-1} = \dfrac{0}{3} = 0$

    - 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   $|1-0| = \left|\dfrac{2}{3} - \dfrac{3}{3}\right| = \dfrac{1}{3}$

‹#›