

# Proximity Measure for Binary Attributes

- A contingency table for binary data

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

- Distance measure for **symmetric binary** variables

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

- Distance measure for **asymmetric binary** variables:

$$d(i, j) = \frac{r + s}{q + r + s}$$

- Jaccard coefficient (**similarity** measure for *asymmetric* binary variables):

$$sim_{Jaccard}(i, j) = \frac{q}{q + r + s}$$

- Note: Jaccard coefficient is the same as

(a concept discussed in Pattern Discovery)

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

# Example: Dissimilarity between Asymmetric Binary Variables

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	Male	Y <sup>yes/no</sup>	N <sup>egative</sup>	P	N	N	N
Mary	Female	Y	N	P	N	P	N
Jim	M	Y	P <sup>ositive</sup>	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:  $d(i, j) = \frac{r + s}{q + r + s}$

$$d(jack, mary) = \frac{0 + 1}{2 + 0 + 1} = 0.33 \frac{1}{3}$$

$$d(jack, jim) = \frac{1 + 1}{1 + 1 + 1} = 0.67 \frac{2}{3}$$

$$d(jim, mary) = \frac{1 + 2}{1 + 1 + 2} = 0.75 \frac{3}{4}$$

		Mary		
		1	0	$\Sigma_{\text{row}}$
Jack	1	2	0	2
	0	1	3	4
	$\Sigma_{\text{col}}$	3	3	6

		Jim		
		1	0	$\Sigma_{\text{row}}$
Jack	1	1	1	2
	0	1	3	4
	$\Sigma_{\text{col}}$	2	4	6

		Mary		
		1	0	$\Sigma_{\text{row}}$
Jim	1	1	1	2
	0	2	2	4
	$\Sigma_{\text{col}}$	3	3	6

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	Y	P	N	N	N	N

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

Symmetric binary

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

$$= \frac{1 + 1}{2 + 1 + 1 + 3}$$

$$= \frac{2}{7}$$

Mary

Jack

	1	0	Sum
1	" 2 <sub>q</sub>	' 1 <sub>r</sub>	3
0	' 1 <sub>s</sub>	"" 3 <sub>t</sub>	4
Sum	3	4	7

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	Y	N	P	N	P	N
Jim	M 1	Y 1	P 1	N 0	N 0	N 0	N 0

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

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

Symmetric binary min 162 = 2  
7

Jack

Jim

	1	0	sum
1	$q$	$r$	
0	$s$	$t$	
sum			

# Proximity Measure for Categorical Attributes

เป็นชื่อที่ใกล้เคียง

- Categorical data, also called nominal attributes
  - Example: Color (red, yellow, blue, green), profession, etc.
- Method 1: Simple matching
  - $m$ : # of matches,  $p$ : total # of variables

$p$  = Attributes ทั้งหมด

$m$  = จำนวนที่เหมือนกัน

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

จำนวนที่เหมือนกัน

จำนวนทั้งหมด

- Method 2: Use a large number of binary attributes
  - Creating a new binary attribute for each of the  $M$  nominal states

Dummy

## Method 2

Maximum Categories in the binary

**๑**      **วิธี**

r	400.
r	07.
g	400.

$$r, g, b$$

עס זאגט, און  
אבר, ערד

	2 R	2 G	2 B	2 Y	2 O	2 I	2 V
1	1	0	0	0	1	0	0
2	1	0	0	0	0	1	0
3	0	1	0	0	1	0	0

၄-ပေးကံ၊ ၅-အကံ၊ ၁ နှစ် ၃ ကြိမ် ပုံစံ (Binary)

Ans  $\frac{2}{7}$

# Ordinal Variables

- An ordinal variable can be discrete or continuous
- Order is important, e.g., rank (e.g., freshman, sophomore, junior, senior)  
 $\{1, 2, 3, 4\}$
- Can be treated like interval-scaled

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

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

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

$\frac{1-1}{4-1} = \frac{0}{3} = 0$ 
 $\frac{2-1}{4-1} = \frac{1}{3}$

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

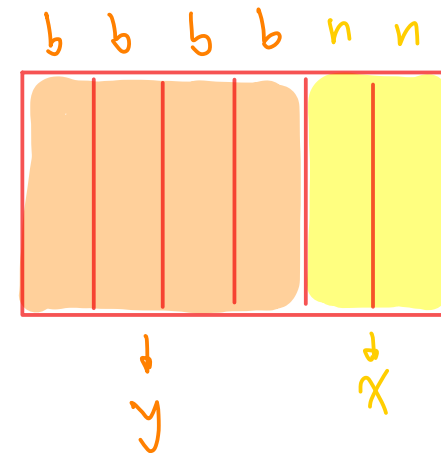
□ Then distance:  $d(\text{freshman}, \text{senior}) = 1$ ,  $d(\text{junior}, \text{senior}) = 1/3$

□ Compute the dissimilarity using methods for interval-scaled variables

# Attributes of Mixed Type

- A dataset may contain all attribute types
  - Nominal, symmetric binary, asymmetric binary, numeric, and ordinal
- One may use a weighted formula to combine their effects:

$$d(i, j) = \frac{\sum_{f=1}^p w_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^p w_{ij}^{(f)}} \rightarrow y\left(\frac{4}{6}\right), x\left(\frac{2}{6}\right)$$



- If  $f$  is numeric: Use the **normalized** distance
- If  $f$  is binary or nominal:  $d_{ij}^{(f)} = 0$  if  $x_{if} = x_{jf}$ ; or  $d_{ij}^{(f)} = 1$  otherwise
- If  $f$  is ordinal

- Compute ranks  $z_{if}$  (where  $z_{if} = \frac{r_{if} - 1}{M_f - 1}$ )

- Treat  $z_{if}$  as interval-scaled



# Cosine Similarity of Two Vectors

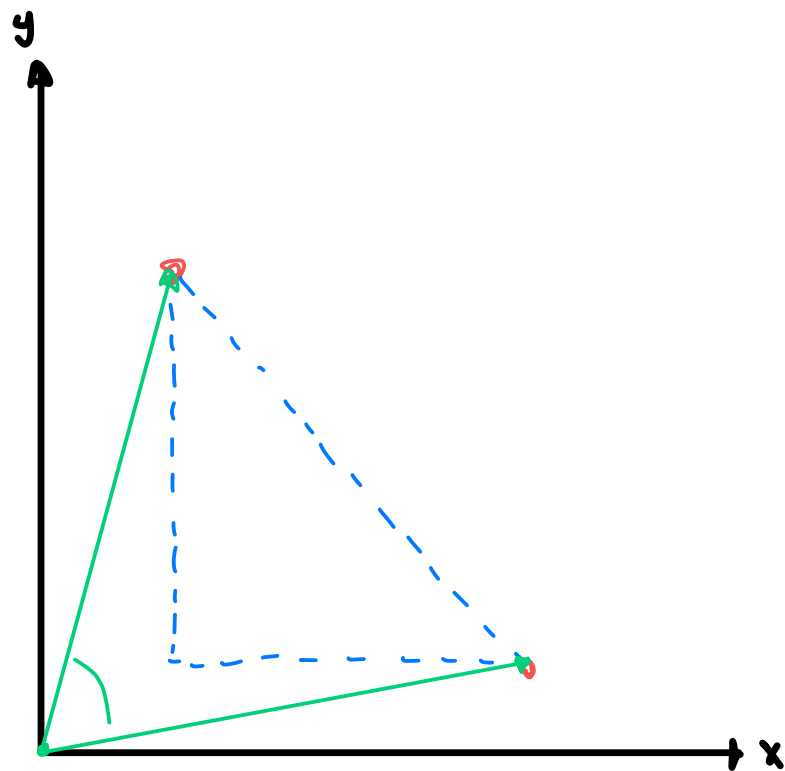
- A **document** can be represented by a bag of terms or a long vector, with each attribute recording the *frequency* of a particular term (such as word, keyword, or phrase) in the document

Document	team	coach	hockey	baseball	soccer	penalty	score	win	loss	season
Document1	5	0	3	0	2	0	0	2	0	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	0
Document4	0	1	0	0	1	2	2	0	3	0

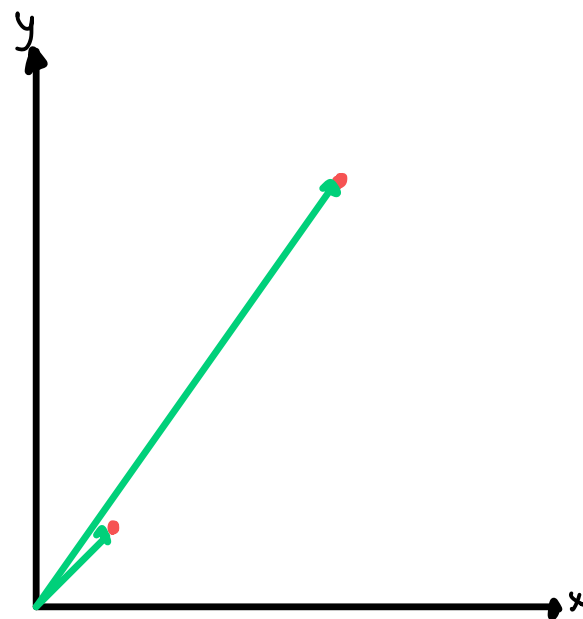
- Other **vector** objects: Gene features in micro-arrays
- Applications: Information retrieval, biologic taxonomy, gene feature mapping, etc.
- Cosine measure: If  $d_1$  and  $d_2$  are two vectors (e.g., term-frequency vectors), then

$$\cos(d_1, d_2) = \frac{d_1 \bullet d_2}{\|d_1\| \times \|d_2\|}$$

where  $\bullet$  indicates vector dot product,  $\|d\|$ : the length of vector  $d$



သော  $\rho_1 + \rho_2$  ပုံသဏ္ဌာန် အောက်ပါ



သော  $\rho_1 + \rho_2$  ပုံသဏ္ဌာန် အောက်ပါ