

Proximity Measure for Binary Attributes

- A contingency table for binary data

		Object <i>j</i>		
		1	0	sum
Object <i>i</i>	1	<i>q</i>	<i>r</i>	<i>q</i> + <i>r</i>
	0	<i>s</i>	<i>t</i>	<i>s</i> + <i>t</i>
	sum	<i>q</i> + <i>s</i>	<i>r</i> + <i>t</i>	<i>p</i>

đi
0 ñú 1

- 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):

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

- Note: Jaccard coefficient is the same as

(a concept discussed in Pattern Discovery)

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

Example: Dissimilarity between Asymmetric Binary Variables

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P = Positive	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N = Negative	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: $d(i, j) = \frac{r + s}{q + r + s}$

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

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

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

		Mary		
		1	0	Σ_{row}
Jack	1	2	0	2
	0	1	3	4
Σ_{col}		3	3	6

		Jim		
		1	0	Σ_{row}
Jack	1	1	1	2
	0	1	3	4
Σ_{col}		2	4	6

		Mary		
		1	0	Σ_{row}
Jim	1	1	1	2
	0	2	2	4
Σ_{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 - 1	Y - 1	P - 1	N - 0	N - 0	N - 0	N - 0

mary

	1	0	SUM
Jack	1	2 9	1 8
	0	1 5	3 4
SUM	3	4	

~~※※※※※, မေး၏~~

Symmetric

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

$$= \frac{1+1}{2+1+1+3} = \frac{2}{7} \neq$$

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

Jack နှင့် Jim

မြတ်နောက် $\frac{2}{7}$ ရမည်။

Symmetric ဆုံးဖို့ပြုလိုက်မည်။

Proximity Measure for Categorical Attributes

- Categorical data, also called nominal attributes

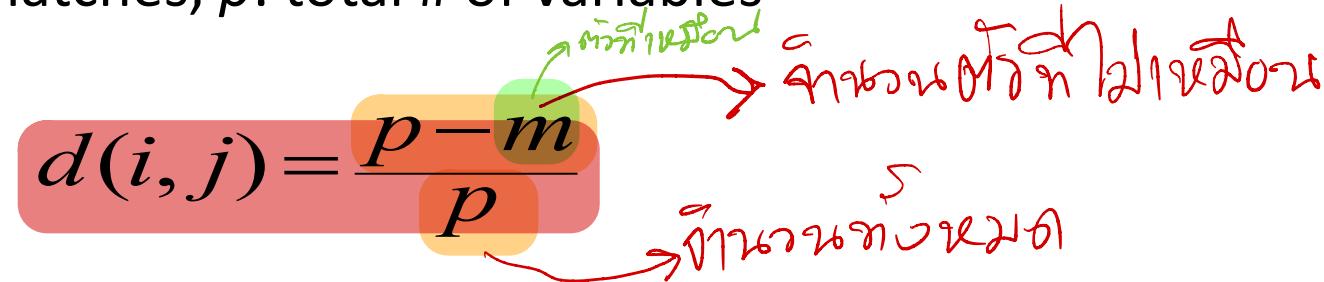
- Example: Color (red, yellow, blue, green), profession, etc.

ස්ථාන තැන්ව සේම මෙය නීතිවාගිකුත් වේ Categorical

- Method 1: Simple matching

- m : # of matches, p : total # of variables

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



- Method 2: Use a large number of binary attributes

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

1

0 1 0

నేడు dummy లో OneHot Encoder

	0 1 0
r	0 0 1
r	0 0 1
g	0 1 0

r, g, b 0, 1, 0, 1, Grab

0 R	0 G	0 B	0	1	0	0	0
1	1	0	0	1	0	0	0
1	0	0	0	0	1	0	0

Ordinal Variables

ຕາງປະກາດ

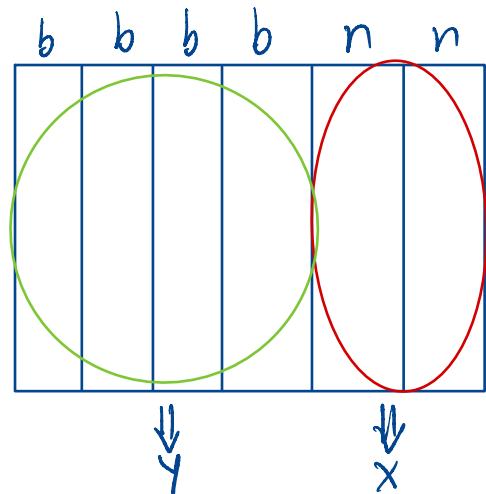
- An ordinal variable can be discrete or continuous
- Order is important, e.g., rank (e.g., freshman, sophomore, junior, senior)
 - 1 2 3 4
- Can be treated like interval-scaled
 - Replace *an ordinal variable value* by its rank:
 $r_{if} \in \{1, \dots, M_f\}$
 - Map the range of each variable onto [0, 1] by replacing *i*-th object in the *f*-th variable by
- Example: freshman: 0; sophomore: 1/3; junior: 2/3; senior 1
 - $$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$
 - freshman
 $\frac{1-1}{4-1} = \frac{0}{3} = 0$
- Then distance: $d(\text{freshman}, \text{senior}) = 1$, $d(\text{junior}, \text{senior}) = 1/3 \Rightarrow$ သົດລະການ
absolute (—)
- 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
ภาษาพื้นฐาน attribute คือ
- One may use a **weighted** formula to combine their effects:

$$d(i, j) = \frac{\sum_{f=1}^p w_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^p w_{ij}^{(f)}}$$

for Distance กับชุดค่าของ Column
คือ ตัวอย่างที่ไม่ต้องคำนึง Column
(weigh)



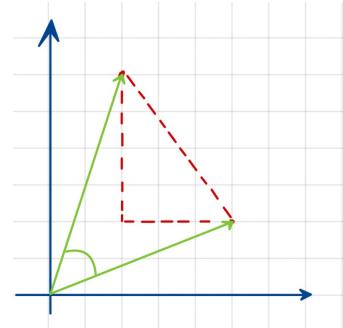
- If f is numeric: Use the normalized distance
- If f is binary or nominal: $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$; or $d_{ij}^{(f)} = 1$ otherwise
- If f is ordinal
 - Compute ranks z_{if} (where $z_{if} = \frac{r_{if} - 1}{M_f - 1}$)
 - Treat z_{if} as interval-scaled

Cosine Similarity of Two Vectors

សំណិតឱ្យទាក់ទងការបង្ហាញសាស្ត្រ នូវករណ៍ពេលការបង្ហាញសាស្ត្រ normal និងការបង្ហាញសាស្ត្រដែលមានភាពខ្លាំងខ្លះ, ហើយនូវករណ៍ពេលការបង្ហាញសាស្ត្រ

- 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

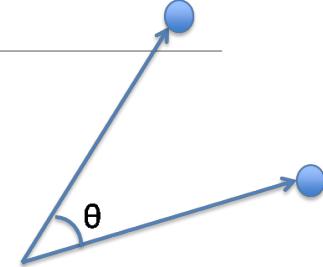
ឯកសារ
គោលការណ៍ដែលបានបង្ហាញសាស្ត្រ

Example: Calculating Cosine Similarity

- Calculating Cosine Similarity:

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

$$sim(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$



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

- Ex: Find the **similarity** between documents 1 and 2.

$$d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0) \quad d_2 = (3, 0, 2, 0, 1, 1, 1, 0, 1, 0)$$

- First, calculate vector dot product

$$d_1 \bullet d_2 = 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 1 + 2 \times 1 + 0 \times 0 + 0 \times 1 = 25$$

- Then, calculate $\|d_1\|$ and $\|d_2\|$

$$\|d_1\| = \sqrt{5 \times 5 + 0 \times 0 + 3 \times 3 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 0 \times 0} = 6.481$$

$$\|d_2\| = \sqrt{3 \times 3 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 1 \times 1 + 1 \times 1 + 0 \times 0 + 1 \times 1 + 0 \times 0 + 1 \times 1} = 4.12$$

- Calculate cosine similarity: $\cos(d_1, d_2) = 25 / (6.481 \times 4.12) = 0.94$