

Proximity Measure for Binary Attributes

ค่า 0 กับ 1

- A contingency table for binary data

		Object j		
		1	0	sum
Object i	1	จำนวนที่ 1 ตรงกัน q หรือ 100	r 0 กับ 1	$q + r$
	0	1 กับ 0 s 100	t 0 กับ 0	$s + t$
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	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

Proximity Measure for Categorical Attributes

- Categorical data, also called nominal attributes

- Example: ชื่อ ex. สีเหลือง หรือ อาชีพ หรืออะไรก็ได้ที่เข้า categorical Color (red, yellow, blue, green), profession, etc.

- Method 1: Simple matching

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

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

ตัวที่เหมือน
จำนวนที่ไม่เหมือน
จำนวนทั้งหมด

- Method 2: Use a large number of binary attributes

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

Ordinal Variables

ตัวแปรอันดับ

- ❑ An ordinal variable can be discrete or continuous
- ❑ Order is important, e.g., rank (e.g., freshman¹, sophomore², junior³, senior⁴)
- ❑ Can be treated like interval-scaled
 - ❑ Replace *an ordinal variable value* by its rank: $r_{if} \in \{1, \dots, M_f\}$
ลำดับที่เท่าไรของในแต่ละชื่อหลังจากนำมาเรียงกัน
ตัวสุดท้ายหลังจากเรียงกัน
 - ❑ Map the range of each variable onto [0, 1] by replacing i -th object in the f -th variable by $z_{if} = \frac{r_{if} - 1}{M_f - 1}$
สูตรหารระยะห่างระหว่างจุด
2 จุด หารระยะห่างด้วยการแทนค่า z
 - ❑ Example: freshman: 0; sophomore: 1/3; junior: 2/3; senior 1
 - ❑ Then distance: $d(\text{freshman}, \text{senior}) = 1$, $d(\text{junior}, \text{senior}) = 1/3$ หาได้โดยการใส่ absolute (|I|)
 - ❑ 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:

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

เอา distance ทั้งหมดคำนวณ และ ถ่วงน้ำหนัก (weigh) ด้วย column

- 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

- นับคำศัพท์จากบทความ สเกลจะเท่าๆกัน , ไม่สามารถ normal ได้ถ้ามี document ที่เป็นหนังสือเพิ่มเข้ามา เพราะ จำนวนคำเยอะ , ทำให้สเกลไม่เท่ากันได้
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

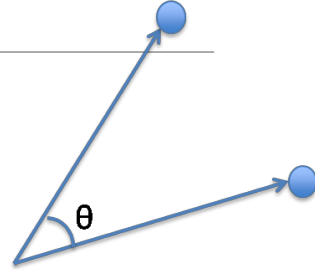
- Other vector objects: Gene features in micro-arrays
- Applications: Information retrieval, biologic taxonomy, gene feature mapping, etc.
- Cosine measure: If d_1 and d_2 are two vectors (e.g., term-frequency vectors), then

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

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

Example: Calculating Cosine Similarity

□ Calculating Cosine Similarity: $\cos(d_1, d_2) = \frac{d_1 \bullet d_2}{\|d_1\| \times \|d_2\|}$ $\text{sim}(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$



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

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

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

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