

12/01/21    dec-07

- ✓ 5. Proximity measures for ordinal attributes
- ✓ 6. Proximity measures for mixed-type "
- ✓ 7. Cosine similarity

5) Ordinal attributes :

Step 1) Let ordinal attributes be ranked as :

< fair, good, excellent >

∴  $m_g$  represents total no. of ordered states  
Here it is 03.

Step 2 : Replace each ordinal data by a rank .  
< fair : 01    good : 02    excellent : 03 >

Step 3 : Normalize the rankings

$$Z_{ij} = \frac{x_{ij} - 1}{m_g - 1}$$

$$\text{fair (1)} = \frac{1-1}{3-1} = 0$$

$$\text{Good (2)} = \frac{2-1}{3-1} = 0.5$$

$$\text{Excellent (3)} = \frac{3-1}{3-1} = 1$$

Eg:

obj-id	Test
1	Excellent (01)
2	fair (0)
3	good (0.5)
4	excellent (1)

$$\begin{pmatrix} 1 & 2 & 3 & 4 \\ 1 & 0 & & \\ 2 & 1 & 0 & \\ 3 & 0.5 & 0.5 & 0 \\ 4 & 0 & 1 & 0.5 & 0 \end{pmatrix}$$

$$d(2,1) = |0-1| = 1$$

$$d(3,1) = |0.5-1| = 0.5$$

$$d(3,2) = 0.5 - 0 = 0.5$$

$$d(4,1) = 0$$

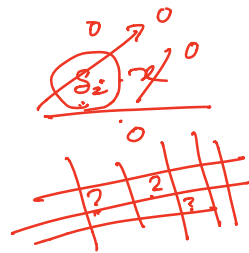
$$d(4,2) = 1$$

$$d(4,3) = 0.5$$

$$d(x,y) = |x_1 - y_1| + |x_2 - y_2| + \dots + |x_n - y_n|$$

6) Proximity measure for mixed type attributes:

Obj id	Test-1 Nominal	Test-2 ordinal	Test-3 Numerical
1	code A	Excellent	45
2	code B	fair	22
3	code C	good	64
4	code A	excellent	28



1. dissimilarity matrix for nominal :  $d(i,j) = \frac{P-m}{P}$

	1	2	3	4
1	0			
2	1	0		
3	1	1	0	
4	0	1	1	0

where  $P = \# \text{ attributes}$   
 $m = \# \text{ match b/w } i \text{ \& } j$

2. dissimilarity matrix for ordinal attribute :

	1	2	3	4
1	0			
2	1	0		
3	0.5	0.5	0	
4	0	1	0.5	0

$$Z_{ij} = \frac{x_{ij} - 1}{m_j - 1}$$

3. dissimilarity matrix for numerical data :

To normalize, the data is mapped in the range of  $[0,1]$ .

$$d(i,j) = \frac{|x_{if} - x_{jf}|}{\max - \min}$$

id	
1	→ 45
2	→ 22
3	→ 64
4	→ 28

	1	2	3	4
1	0			
2	0.55	0		
3	0.45	0.1	0	
4	0.4	0.14	0.8	0

$$d(2,1) = \frac{|45 - 22|}{64 - 22} = 0.55$$

$$d(3,1) = \frac{|64 - 45|}{64 - 22} = 0.45$$

$$d(3,2) = \frac{|64 - 22|}{64 - 22} = 0.1$$

4) formula for the mixed attribute:

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

where  $\delta_{ij} = 0$  if  $x_{ij}$  or  $x_{jf}$  is missing or  $x_{ij} = x_{jf} = 0$  and  $f$  is assy. binary  
else  $\delta_{ij} = 1$

$$\begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{pmatrix} 0 & & & \\ 0.85 & 0 & & \\ 0.65 & 0.83 & 0 & \\ 0.13 & 0.71 & 0.79 & 0 \end{pmatrix} \end{matrix}$$

$$d(2,1) = \frac{(1 \times 1) + (1 \times 1) + (1 \times 0.55)}{1+1+1} = \frac{2.55}{3} = 0.85$$

7) Cosine Similarity: & (Cosine distance):

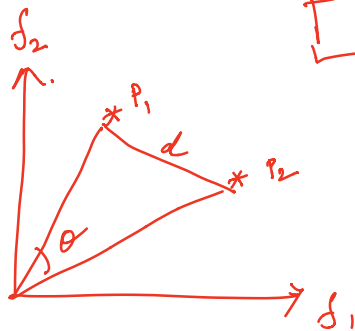
→ widely used in recommendation system (NLP)

$\theta$   $P_1$   $P_2$

dis ↑ sim ↓

dis ↓ sim ↑

$$1 - \cos \theta = \text{cos-distance}$$



if  $\theta = 45^\circ$

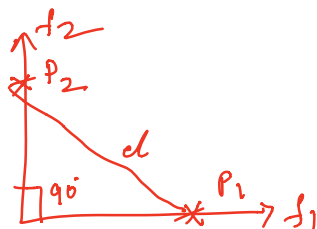
$$\therefore \cos \theta = \cos 45^\circ = 0.53$$

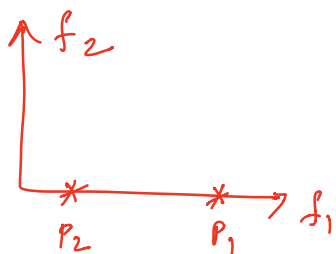
(This means  $P_1$  &  $P_2$  are

53% similar based on angle)

$$\begin{aligned} \cos\text{-dis} &= 1 - \cos \theta \\ &= 1 - 0.53 \\ &= \underline{\underline{0.47}} \end{aligned}$$

$$\cos(90^\circ) = 0$$





$$\cos(0^\circ) = 1$$

2 points are similar as the angle is  $0^\circ$  and similarity is 100%.

### Term - frequency vector :

- Each document is <sup>an</sup> object represented by a term-frequency vector.
- Term frequency are very long and sparse (ie too many 0 values).
- Cosine similarity is a measure of similarity that can be used to compare documents or give a ranking of documents w.r.t a given vector of query words.

$$\text{sim fun: } \boxed{\text{sim}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}}$$

where  $\|x\|$  is the Euclidean norm of vector  $(x_1, x_2, \dots, x_p)$  defined as  $\|x\| = \sqrt{x_1^2 + x_2^2 + \dots + x_p^2}$ . Conceptually, it is the length of the vector.

- A cosine value 0 means the 2 vectors are  $90^\circ$  to each other (orthogonal) hence don't match.
- The closer the cosine value is to 1, the smaller is the angle and greater is the match b/w 2 vectors.

Eg: Term frequency vector :

doc	team	coach	hockey	basket ball	hockey	penalty	score	win	loss	season
1	5	0	3	0	2	0	0	2	0	0
2	3	0	2	0	1	1	0	1	0	1
3	0	7	0	2	1	0	0	3	0	0
4	0	1	0	0	1	2	2	0	3	0

Let  $x = \text{doc1}$  &  $y = \text{doc2}$

find similarity b/w  $x$  &  $y$ .

Soln:  $\text{sim}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$

$$x = \langle 5 \quad 0 \quad 3 \quad 0 \quad 2 \quad 0 \quad 0 \quad 2 \quad 0 \quad 0 \rangle$$

$$y = \langle 3 \quad 0 \quad 2 \quad 0 \quad 1 \quad 1 \quad 0 \quad 1 \quad 0 \quad 1 \rangle$$

$$\begin{aligned} x \cdot y &= (5 \times 3) + (0) + (3 \times 2) + (0) + (2 \times 1) + (0) + (0) + (2 \times 1) + (0) + (0) \\ &+ (2 \times 1) + 0 + 0 = 25 \end{aligned}$$

$$\|x\| = \sqrt{5^2 + 0^2 + 3^2 + 0^2 + 2^2 + 2^2} = 6.48$$

$$\|y\| = \sqrt{3^2 + 0^2 + 2^2 + 0^2 + 1 + 1 + 1 + 1} = 4.12$$

$$\text{sim}(x, y) = \frac{25}{6.48 \times 4.12} = \underline{\underline{0.94}}$$

$\therefore$  2 documents  $x$  &  $y$  are 94% similar