

08/01/21

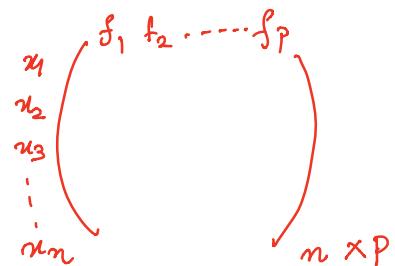
Dec-06

Topics to be discussed :

1. data matrix vs. dissimilarity matrix
2. proximity measure for nominal attribute
3. " " " binary "
4. " " " numeric "
5. " " " ordinal "
6. " " " mixed "
7. cosine similarity

data matrix :

2-mode  
matrix

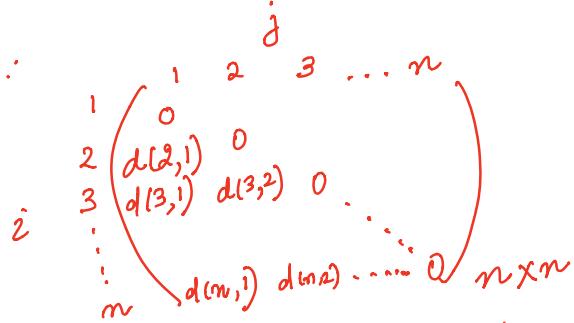


$n \rightarrow$  no. of objects

$p \rightarrow$  " " attributes

(age, ht, wt,  
sal, gender etc.)

dissimilarity :



→ In general,  $d(i,j)$  is a non-negative no. that is close to 0 when objects  $i$  and  $j$  are highly similar or near to each other, and become larger when they differ.

→  $d(i,i) = 0$  (difference between object and itself is 0)

→  $d(i,j) = d(j,i)$  (symmetric)

$$\boxed{\text{sim}(i,j) = 1 - d(i,j)} \quad (\text{for nominal data})$$

2) Proximity measure for nominal attribute :

Nominal attribute :

Color code

Blue  
red  
green  
Yellow

Grading

A  
B  
C  
A

dissimilarity b/w  $i \neq j$ :

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

where  $P = \text{no. of attributes}$

$m = \text{no. of matches}$   
between  $i \neq j$ .

ID	Test Res	$f_1 (P=1)$
1	code A	
2	code B	
3	code C	
4	code A	

$$d(2,1) = \frac{1-0}{1} = 1$$

$$\begin{matrix} & 1 & 2 & 3 & 4 \\ 1 & 0 & & & \\ 2 & d(2,1) & 0 & & \\ 3 & d(3,1) & d(3,2) & 0 & \\ 4 & d(4,1) & d(4,2) & d(4,3) & 0 \end{matrix} \quad 4 \times 4$$

$$d(3,1) = \frac{1-0}{1} = 1$$

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

$$d(3,2) = \frac{1-0}{1} = 1$$

$$d(4,1) = \frac{1-1}{1} = 0$$

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

$$d(4,3) = 1$$

2) Proximity measure for binary attribute:

Sym. binary:  
(Gender: M/F)

assym. binary  
(Medical Test  
Outcome)

$$d(i,j) = \frac{n+s}{n+n+k+l}$$

dissimilarity

sym.

binary

assym.  
binary

$$d(i,j) = \frac{n+s}{n+n+k+l}$$

similarity

sym.

binary

assym.  
binary

\*

Simple matching coefficient

$$\text{Sim}(i,j) = \frac{M_{11} + M_{00}}{M_{11} + M_{10} + M_{01} + M_{00}}$$

\* Jaccard Coefficient

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

$$= \frac{M_{11}}{M_{11} + M_{10} + M_{01}} = 1 - d(i,j)$$

Contingency Table :

		(j)	
		1	0
(i)	1	q "	r "
	0	s "	t "

q : # attributes that equal 01 for obj i & j

r : # " " " " " 01 n & 0 for j

s : # " " " " " 0 " i & 1 " "

t : # " " " " " 0 " " " 0 " j

$$\boxed{\text{Total no. of attributes} = p = q + r + s + t}$$

Eg :

	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>	T <sub>4</sub>	T <sub>5</sub>	T <sub>6</sub>
Jack	1	0	1	0	1	0
Jim	1	0	1	0	0	0
Mary	1	1	0	0	1	0

Care!: Asymm. Binary :  $d(i,j) = \frac{M_{10} + M_{01}}{M_{11} + M_{10} + M_{01}}$

		Jack	Jim	Mary
		0	0.33	0
Jack	0			
Jim	0.33			

$$\begin{array}{c}
 \text{Mary} \backslash 0.66 & 0.75 & 0 / \\
 \text{Tim} \quad \text{Jack} \quad \text{John} \\
 \end{array}$$

$$d(\text{Tim}, \text{Jack}) = \frac{1+0}{2+1+0} = \frac{1}{3} = 0.33$$

$$d(\text{Mary}, \text{Jack}) = 0.66$$

$$d(\text{Mary}, \text{Tim}) = 0.75$$

Case 2: for symm. binary

	Jack	Tim	Mary
Jack	0		
Tim	$\frac{1}{6}$	0	
Mary	$\frac{1}{3}$	$\frac{1}{2}$	0

4) Proximity measure for numeric attribute:

3 most commonly used measures are:

1. Minkowski distance
2. Euclidean distance
3. Manhattan distance
4. Supremum distance

1. Minkowski distance:

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

where  $h$  is a real no. s.t.  $h \geq 1$ . — (1)

2. Euclidean distance:

when  $h=2$ , Eqn(1) becomes  $L_2$  norm i.e. Euclidean distance.

$$d(i, j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2} — (2)$$

where  $i \neq j$  are 2 objects having "p" no. of attributes.

3. Manhattan distance : / or (city block distance)

In eqn. (1), when  $h=1$ , L1 norm or Manhattan distance is derived.

$$d(i,j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}| \quad \text{--- (3)}$$

4) Supremum distance:

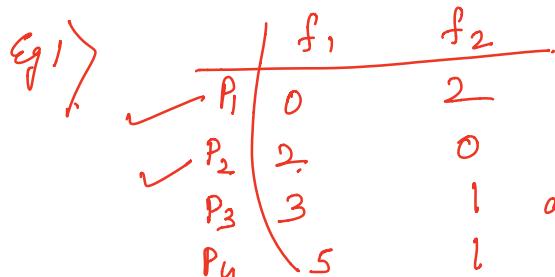
This is also referred as  $L_{\max}$ ,  $L_{\infty}$  norm and Chebyshew distance which is a generalization of Minkowski distance for  $h \rightarrow \infty$ .

→ To compute it, find the attribute of that gives the maximum distance in values b/w 2 objects.

$$d(2,j) = \lim_{h \rightarrow \infty} \left( \sum_{f=1}^P |x_{if} - x_{jf}|^h \right)^{1/h}$$

$$= \max_f |x_{if} - x_{jf}|$$

$L^{\infty}$  norm is also known as uniform norm.

Eg1) 

;  $\Rightarrow$  Euclidean :

$d(P_2, P_1) =$	$\sqrt{(2-0)^2 + (0-2)^2}$
	$= \sqrt{4+4} = \sqrt{8}$

$P_1 \begin{pmatrix} 0 & 2 \\ 2 & 0 \\ 3 & 1 \\ 5 & 1 \end{pmatrix}, P_2 \begin{pmatrix} 0 & 2.8 \\ 2.8 & 0 \\ 3.2 & 1.4 \\ 5.1 & 3.2 \end{pmatrix}, P_3 \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}, P_4 \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}$

(ii) Manhattan:

$$d(P_2, P_1) = |2-0| + |0-2| \\ = 2+2 \\ = 4$$

$$P_1 \begin{pmatrix} 0 & P_1 & P_2 & P_3 & P_4 \\ P_2 & 4 & 0 & 0 & 0 \\ P_3 & 4 & 2 & 0 & 0 \\ P_4 & 6 & 4 & 2 & 0 \end{pmatrix}$$

(iii) Supremum distance :

$$d(P_2, P_1) = \max_{i \in \delta} \{ |2 - 0|, |0 - 2| \}$$
$$= \max \{ 2, 2 \}$$
$$= 2$$

	$P_1$	$P_2$	$P_3$	$P_4$
$P_1$	0			
$P_2$		0		
$P_3$	3	1	0	
$P_4$	5	3	2	0

.....