

Similarity, Dissimilarity, and Proximity

ความคล้ายคลึง

□ Similarity measure or similarity function

- A real-valued function that quantifies the similarity between two objects
- Measure how two data objects are alike: The higher value, the more alike
- Often falls in the range [0,1]: 0: no similarity; 1: completely similar

ความต่างกัน
ดัชนีความคล้ายคลึง

□ Dissimilarity (or distance) measure

- Numerical measure of how different two data objects are
- In some sense, the inverse of similarity: The lower, the more alike
- Minimum dissimilarity is often 0 (i.e., completely similar) เนื่องจากว่า 0
[เนื่อง, ทำ]
- Range [0, 1] or [0, ∞), depending on the definition

ความต่าง

□ Proximity usually refers to either similarity or dissimilarity

Data Matrix and Dissimilarity Matrix

- Data matrix

- A data matrix of n data points with l dimensions

$$D = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1l} \\ x_{21} & x_{22} & \dots & x_{2l} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nl} \end{pmatrix}$$

- Dissimilarity (distance) matrix

- n data points, but registers only the distance $d(i, j)$ (typically metric)

$$\begin{pmatrix} 0 & & & \\ d(2,1) & 0 & & \\ \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & 0 \end{pmatrix}$$

- Usually symmetric, thus a triangular matrix

- **Distance functions** are usually different for real, boolean, categorical, ordinal, ratio, and vector variables

- Weights can be associated with different variables based on applications and data semantics

Standardizing Numeric Data

- Z-score:

$$z = \frac{x - \mu}{\sigma}$$

- X: raw score to be standardized, μ : mean of the population, σ : standard deviation
- the distance between the raw score and the population mean in units of the standard deviation
- negative when the raw score is below the mean, "+" when above
- An alternative way: Calculate the mean absolute deviation

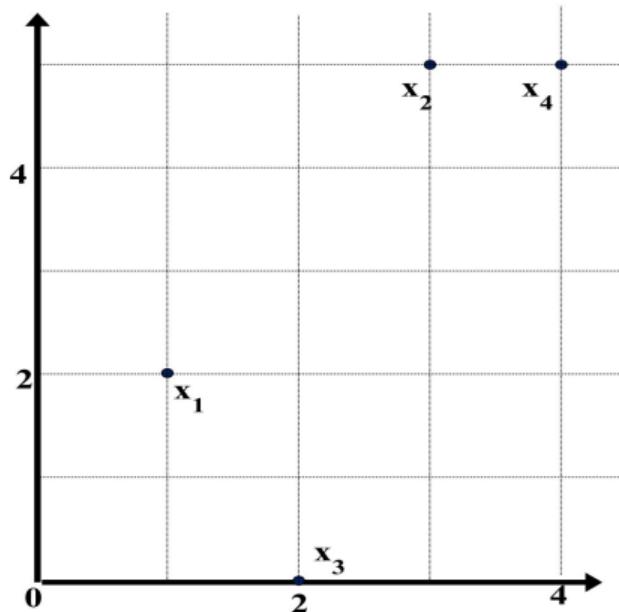
$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

where

$$m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf})$$

- standardized measure (z-score):
$$z_{if} = \frac{x_{if} - m_f}{s_f}$$
- Using mean absolute deviation is more robust than using standard deviation

Example: Data Matrix and Dissimilarity Matrix



Data Matrix

point	attribute1	attribute2
$x1$	1	2
$x2$	3	5
$x3$	2	0
$x4$	4	5

Dissimilarity Matrix (by Euclidean Distance)

	$x1$	$x2$	$x3$	$x4$
$x1$	0			
$x2$	3.61	0		
$x3$	2.24	5.1	0	
$x4$	4.24	1	5.39	0

Distance on Numeric Data: Minkowski Distance

- ❑ **Minkowski distance**: A popular distance measure

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

where $i = (x_{i1}, x_{i2}, \dots, x_{il})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jl})$ are two l -dimensional data objects, and p is the order (the distance so defined is also called L- p norm)

- ❑ Properties
 - ❑ $d(i, j) > 0$ if $i \neq j$, and $d(i, i) = 0$ (Positivity)
 - ❑ $d(i, j) = d(j, i)$ (Symmetry)
 - ❑ $d(i, j) \leq d(i, k) + d(k, j)$ (Triangle Inequality)
- ❑ A distance that satisfies these properties is a **metric**
- ❑ Note: There are nonmetric dissimilarities, e.g., set differences

Special Cases of Minkowski Distance

- $p = 1$: (L_1 norm) Manhattan (or city block) distance ผลรวมของแต่ละแนวแกน
 - E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \cdots + |x_{il} - x_{jl}|$$

- $p = 2$: (L_2 norm) Euclidean distance $\sqrt{\text{Σ[(x_{ij} - x_{kj})}^2]}$

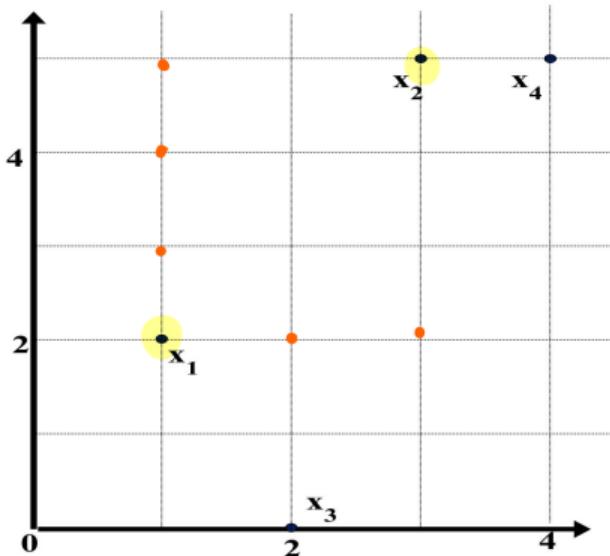
$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \cdots + |x_{il} - x_{jl}|^2}$$

- $p \rightarrow \infty$: (L_{\max} norm, L_{∞} norm) “supremum” distance เคี่ยวก้าว max
 - The maximum difference between any component (attribute) of the vectors

$$d(i, j) = \lim_{p \rightarrow \infty} \sqrt[p]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \cdots + |x_{il} - x_{jl}|^p} = \max_{f=1}^l |x_{if} - x_{jf}|$$

Example: Minkowski Distance at Special Cases

point	attribute 1	attribute 2
x1	1	2
x2	3	5
x3	2	0
x4	4	5



Manhattan (L_1)

L	x1	x2	x3	x4
x1	0			
x2	5	0		
x3	3	6	0	
x4	6	1	7	0

Euclidean (L_2)

L_2	x1	x2	x3	x4
x1	0			
x2	3.61	0		
x3	2.24	5.1	0	
x4	4.24	1	5.39	0

Supremum (L_∞)

L_∞	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0

Proximity Measure for Binary Attributes

- ❑ A contingency table for binary data

		Object <i>j</i>		
		1	0	sum
Object <i>i</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):

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

- ❑ 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	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

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

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

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

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

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

		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

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

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

Jack

	1	0	SUM
1	2 q	1 r	3
0	1 s	3 t	4
SUM	3	4	7

Proximity Measure for Categorical Attributes

- ❑ 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

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

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

Ordinal Variables

- An ordinal variable can be discrete or continuous
- Order is important, e.g., rank (e.g., freshman, sophomore, junior, senior)
- 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
 - Then distance: $d(\text{freshman}, \text{senior}) = 1$, $d(\text{junior}, \text{senior}) = 1/3$
 - Compute the dissimilarity using methods for interval-scaled variables

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$
$$\frac{1-1}{4-1} = \frac{0}{3} = 0$$
$$\frac{2-1}{4-1} = \frac{1}{3}$$