

## Measuring Data Similarity, and Proximity

ข้อมูลที่เป็นตัวเลข สามารถนำมา plot เพื่อดูความห่าง สิ่งที่เราจำเป็นต้องจะทำก่อนจะนำข้อมูลไปประมวลผล คือ ต้องสามารถวัดได้ว่า Data จุดที่ 1 กับ Data จุดที่ 2 มันเหมือน หรือ มันต่างกันอย่างไร

## Similarity, Dissimilarity, and Proximity

- ❑ **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
- ❑ **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)
  - ❑ Range  $[0, 1]$  or  $[0, \infty)$ , depending on the definition
- ❑ **Proximity** usually refers to either similarity or dissimilarity

55

- Similarity measure or similarity function (ความเหมือน)
 

อยากรู้ว่าข้อมูลทั้งสองเหมือนหรือต่างกันอย่างไร จึงทำการสร้างฟังก์ชันขึ้นมาหนึ่งฟังก์ชัน ที่ใส่ Data 2 จุดเข้าไป แล้วฟังก์ชันจะทำการคำนวณว่า สองจุดนี้เหมือนหรือต่างกันอย่างไร โดยผล Output ออกมา จะมีค่าอยู่ระหว่าง 0 – 1 ถ้า Similarity มีค่าเป็น 0 หมายความว่าข้อมูลทั้งสองไม่มีความเหมือนกันเลย, Similarity มีค่าเป็น 1 หมายความว่าข้อมูลทั้งสองมีความเหมือนกัน
- Dissimilarity (or distance) measure (ความไม่เหมือน)
 

วัดว่า ระยะห่างของทั้งสองข้อมูลเป็นเท่าไร ระยะห่างน้อย หมายความว่า ทั้งสองข้อมูลมีความเหมือนกันมาก, ระยะห่างมาก หมายความว่าทั้งสองข้อมูลมีความเหมือนกันน้อย
- Proximity (ความต่าง ระยะห่างระหว่างข้อมูล)

## Data Matrix and Dissimilarity Matrix

### Data matrix

- A data matrix of  $n$  data points with  $l$  dimensions

### Dissimilarity (distance) matrix

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

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

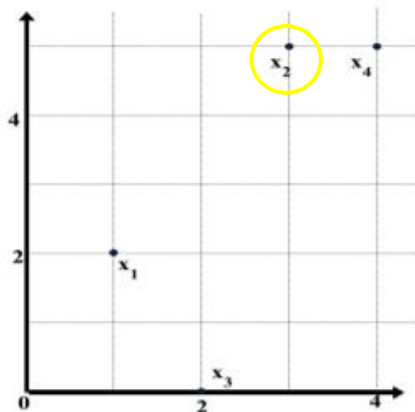
$$\begin{pmatrix} 0 & & & \\ d(2,1) & 0 & & \\ \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & 0 \end{pmatrix}$$

56

แนวนอน เป็นข้อมูลแต่ละจุด, แนวตั้ง เป็นพิวลแต่ละอัน

Distance matrix เป็น matrix ที่ใช้คำนวณก่อนที่จะนำไปประมวลผล ว่า ข้อมูลไหน ห่างจากข้อมูลไหน เท่าไหร่

## 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

58

Dissimilarity Matrix เป็นการบอกว่าข้อมูลแต่ละจุดห่างกันเท่าไหร่ ซึ่งถ้าข้อมูลมี 4 จุด เท่ากับ Distance matrix จะมีขนาดเท่ากับ 4\*4

## Standardizing Numeric Data

- Z-score:  $z = \frac{x - \mu}{\sigma}$ 
  - X: raw score to be standardized,  $\mu$ : mean of the population,  $\sigma$ : 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

57

## 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

59

## Special Cases of Minkowski Distance

□  $p = 1$ : ( $L_1$  norm) **Manhattan (or city block) distance**

□ E.g., the Hamming distance: the number of bits that are different between two binary vectors

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

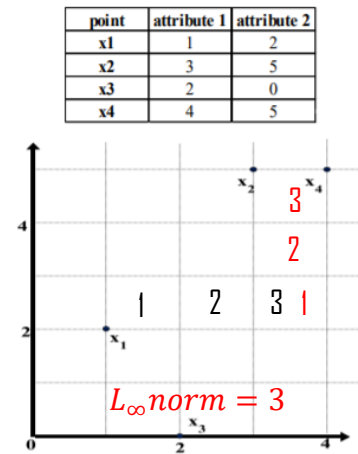
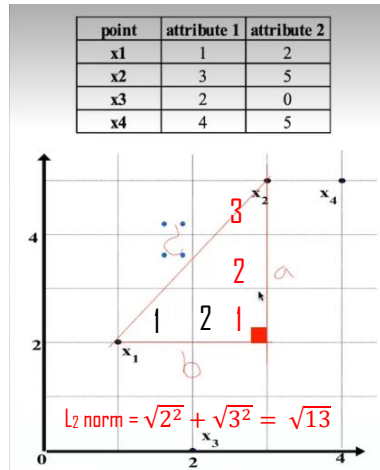
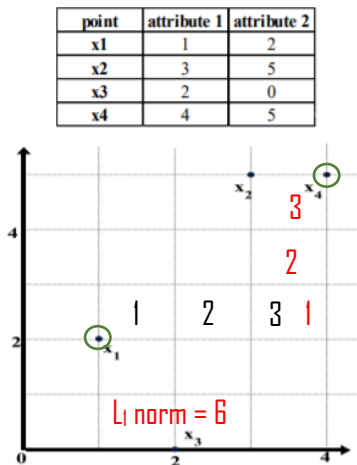
□  $p = 2$ : ( $L_2$  norm) **Euclidean distance**

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

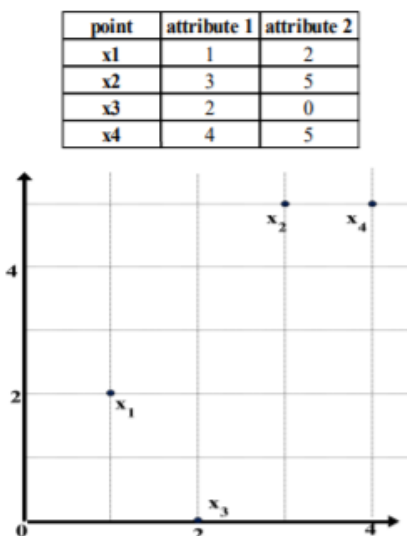
□  $p \rightarrow \infty$ : ( $L_{\max}$  norm,  $L_{\infty}$  norm) **"supremum" distance**

□ The maximum difference between any component (attribute) of the vectors

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



## Example: Minkowski Distance at Special Cases



**Manhattan ( $L_1$ )**

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

**Euclidean ( $L_2$ )**

L2	x1	x2	x3	x4
x1	0			
x2	3.61	0		
x3	2.24	5.1	0	
x4	4.24	1	5.39	0

**Supremum ( $L_{\infty}$ )**

L <sub>∞</sub>	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0

