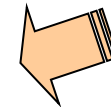


Chapter 2. Getting to Know Your Data

- ❑ Data Objects and Attribute Types
- ❑ Basic Statistical Descriptions of Data
- ❑ Data Visualization
- ❑ Measuring Data Similarity and Dissimilarity
- ❑ Summary

มรที่จจะ Data นี้ไปประมวลผล ก้าวที่จเป็น

เราจะห้ความสามารถให้ Data
จุดที่ 1 กับ Data จุดที่ 2 เปรียบเทียบ
ห้กันห้จวไร ห้ความเหมือนหรือ
ความต่างของ Data ก็จะใช้
Distance หรือ ระยะห่างเป็นมาตรวัด



Similarity, Dissimilarity, and Proximity

ความเหมือน

□ **Similarity** measure or similarity function

↗ สืบพ้องกันเพื่อให้ทราบว่ามีจุดเหมือนกันหรือไม่

□ A real-valued function that quantifies the similarity between two objects

↗ Data ทั้งสองตัว จะเป็นตัวกำหนดความเหมือนว่าจะเท่าไรเหมือนกันหรือไม่

□ Measure how two data objects are alike: The higher value, the more alike

↗ Data ที่มี output จะมีความใกล้เคียง 0-1

□ Often falls in the range $[0,1]$: 0: no similarity; 1: completely similar

ความไม่เหมือน

ใช้จะต่าง

□ **Dissimilarity** (or **distance**) measure

□ Numerical measure of how different two data objects are → ยิ่งไม่เหมือนกันยิ่งต่าง

□ In some sense, the inverse of similarity: The lower, the more alike

□ Minimum dissimilarity is often 0 (i.e., completely similar)

เหมือนกัน

เท่ากัน

↗ ยิ่งต่างยิ่งมีค่ามากขึ้นเรื่อยๆ

□ Range $[0, 1]$ or $[0, \infty)$, depending on the definition

ความต่างหรือจะต่าง (ไม่เหมือน)

□ **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

→ เป็นตัวที่ใช้คำนวณว่า Data ในชุดไหนจาก Data ในชุดไหน

- 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

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

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

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

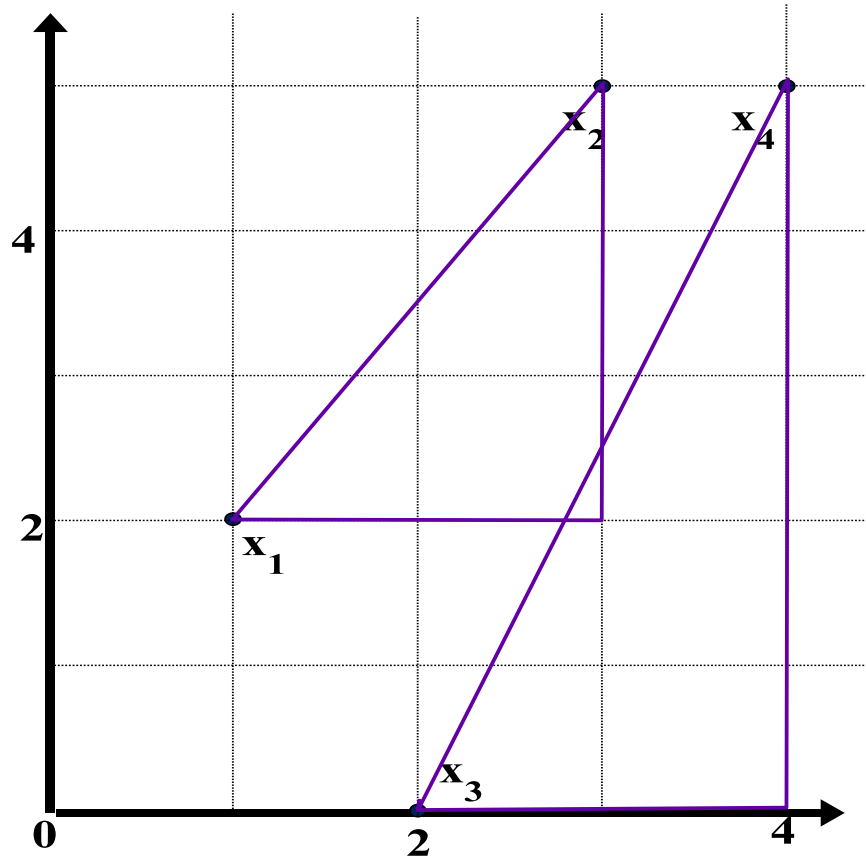
where

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

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

↗ มาตรฐานที่ทำได้ง่ายกว่ามาตรฐาน

Example: Data Matrix and Dissimilarity Matrix



Data Matrix

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

จุด = จุดที่จุดนำวนกันทำใน

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 **มีคุณสมบัติ 3 ข้อ**

- $d(i, j) > 0$ if $i \neq j$, and $d(i, i) = 0$ (**Positivity**) \rightarrow ระยะห่างระหว่างจุด 2 จุด ต้องมากกว่า 0 เสมอ นอกจากริดจุดนั่นเอง

- $d(i, j) = d(j, i)$ (**Symmetry**) x_1 ห่างจาก $x_2 = x_2$ ห่างจาก x_1 เสมอ

- $d(i, j) \leq d(i, k) + d(k, j)$ (**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

เราจะวัดระยะทางในแนวตามแกนอย่างง่าย

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

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

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

□ $p = 2$: (L_2 norm) **Euclidean distance** → ระยะทางที่วัดระหว่างจุดสองจุด โดยใช้ทฤษฎีพีทาโกรัส $c^2 = a^2 + b^2$

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

ระยะทางในระนาบ
จุดแต่ละจุด

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

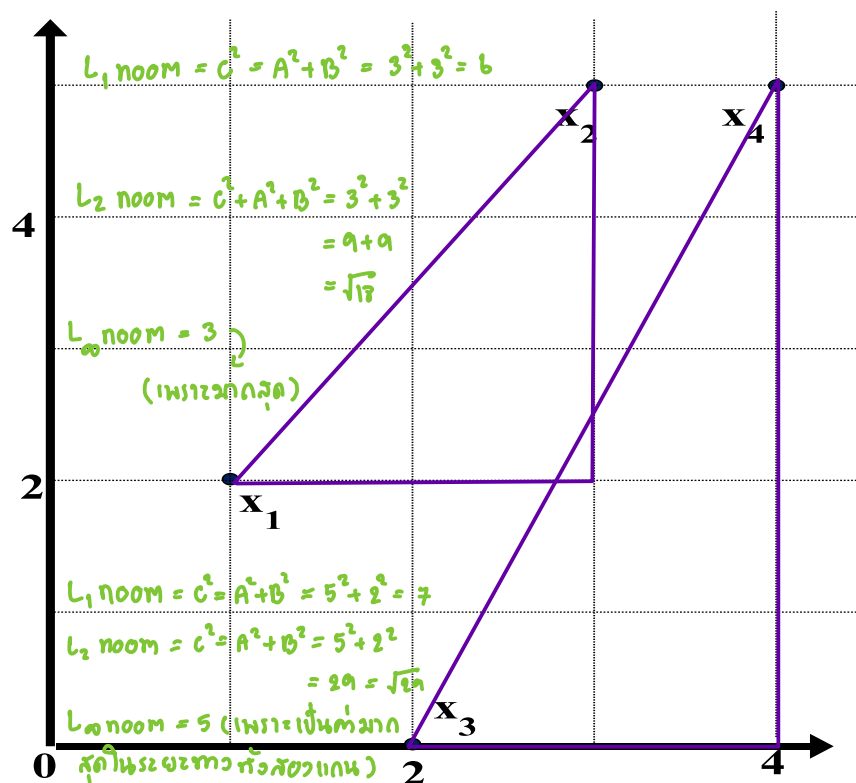
- The maximum difference between any component (attribute) of the vectors

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

สูตรดูหั่นจาก ค่าความจวิ้นว่าค่าที่หั่นสุดฐานนั้น โดยเลือกระยะที่มากที่สุด

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)

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_{∞})

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

$L_{\infty, noom} = \max$ ของระยะทางของทั้งสองแกน (แกนไหนมีค่ามากกว่า)