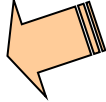


# Chapter 2. Getting to Know Your Data

---

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

# Similarity, Dissimilarity, and Proximity

- **Similarity measure or similarity function**  
*การวัดความเหมือน*  
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  
*จะวัดความเหมือน 1 หมายถึง ที่เราใส่ Data ที่ใกล้เคียงกัน แล้วบอกว่า 2 วัตถุนี้เหมือนกันหรือไม่*  
*Data ที่ output ของการวัดจะเป็นค่าระหว่าง 0-1*  
*- ถ้า similarity เป็น 0 แปลว่าไม่เหมือนกัน*  
*- " " " " เป็น 1 แปลว่าเหมือนกัน*
- **Dissimilarity (or distance) measure**  
*การวัดความต่าง*  
*ระยะห่าง*  
*ที่วัดว่า Data 2 วัตถุต่างกันแค่ไหน*  
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)  
*ความต่าง*  
*ต่างกัน*  
Range  $[0, 1]$  or  $[0, \infty)$ , depending on the definition  
*- ระยะห่างที่ต่างกัน*  
*- ระยะห่างที่เหมือนกัน*
- **Proximity** usually refers to either similarity or dissimilarity  
*ทั้ง similarity หรือ dissimilarity จะวัดความใกล้เคียง*

# Data Matrix and Dissimilarity Matrix

နိမိတ် လိုက်လံသွင်း Data ကို x:၁၀  
လက်တွေ့ feature ကို x:၁၀  
အမှတ်, အမှတ် အညွှန်းကိန်း

## □ 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)



- 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

မှ ပုံသေ

□ Z-score:

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

တိုက်ရိုက် column ၁ မှာရှိတဲ့ အတိုင်း scale ပုံသေ

- 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

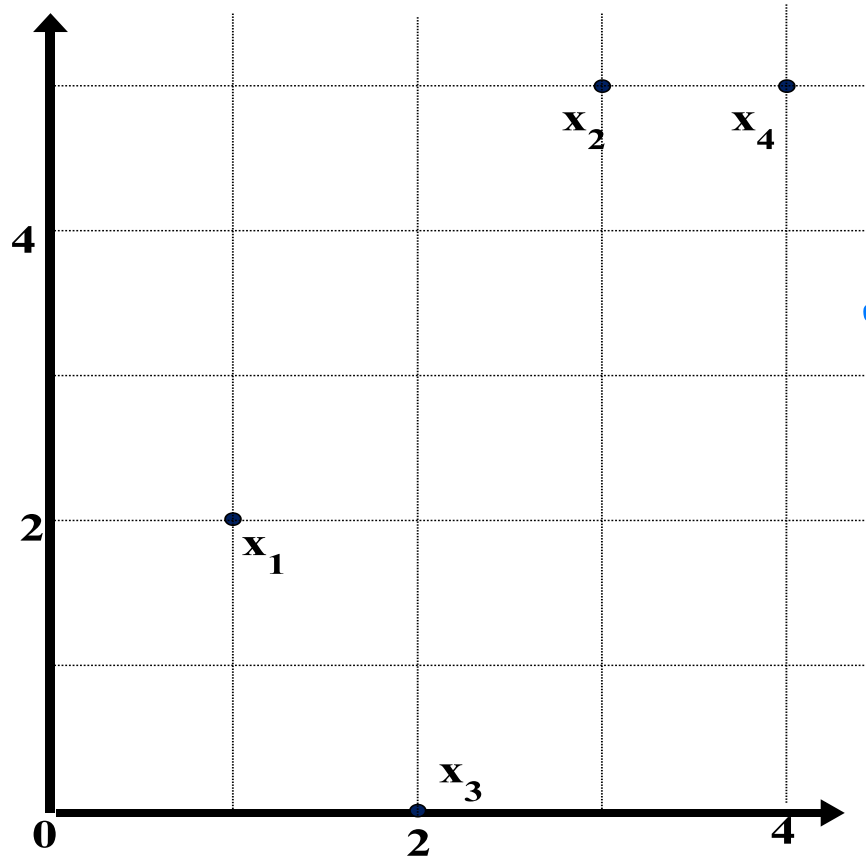
$$m_f = \frac{1}{n} (x_{1f} + x_{2f} + \dots + x_{nf}).$$

□ standardized measure (z-score):

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

□ Using mean absolute deviation is more robust than using standard deviation

# Example: Data Matrix and Dissimilarity Matrix



Data Matrix

2 feature, 2 column

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

matrix of dissimilarity

matrix 4x4

Dissimilarity Matrix (by **Euclidean Distance**)

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

□ 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**

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

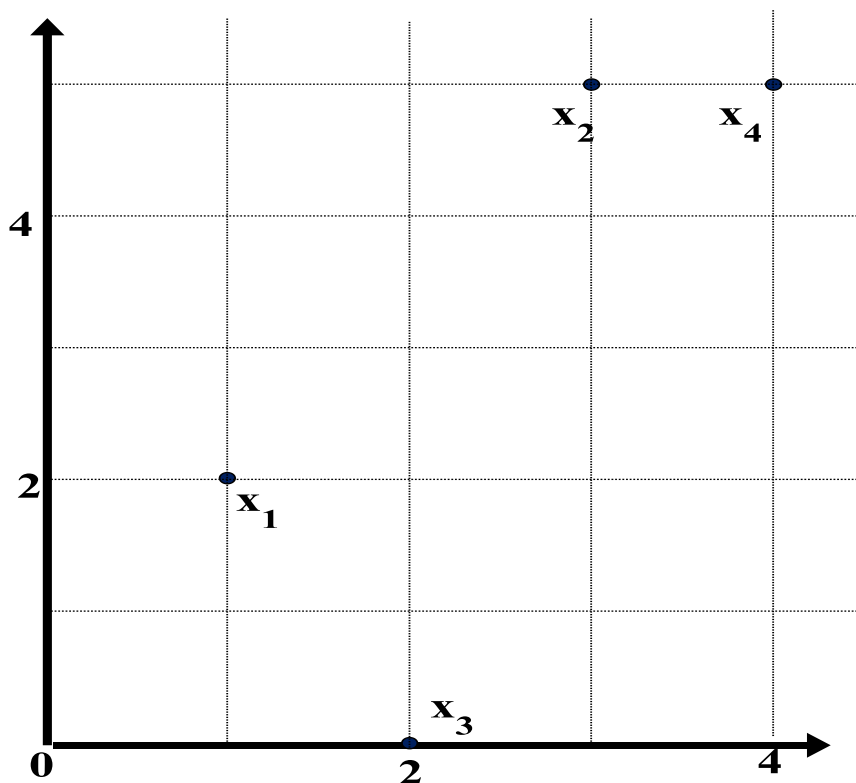
□  $p \rightarrow \infty$ : ( $L_{\max}$  norm,  $L_{\infty}$  norm) **“supremum” distance**

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

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

# 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 ( $L_2$ )

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

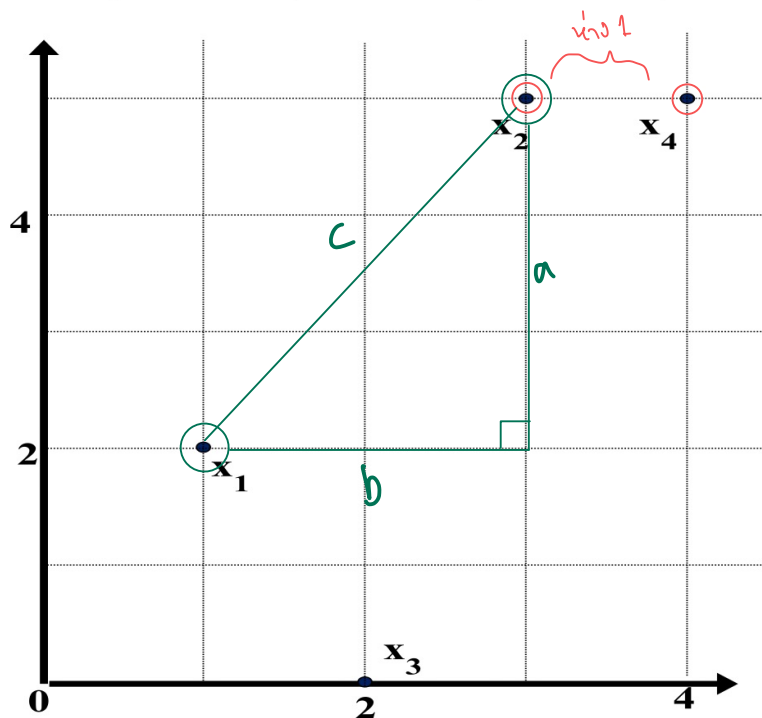
## Supremum ( $L_\infty$ )

$L_\infty$	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0



Example

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



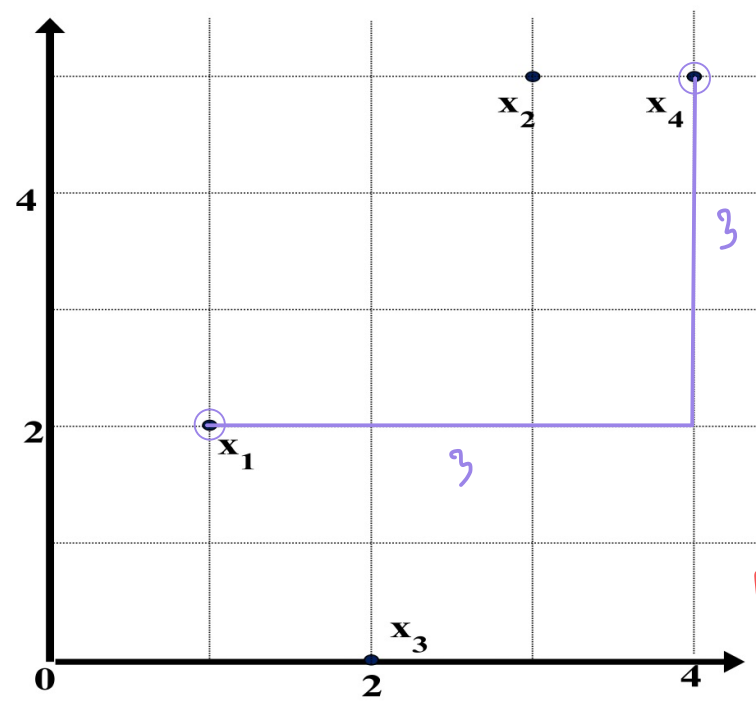
4 จุดที่มีพิกัด  
 $c^2 = a^2 + b^2$   
 $= 3^2 + 2^2$

$= 9 + 4$

$c^2 = 13$

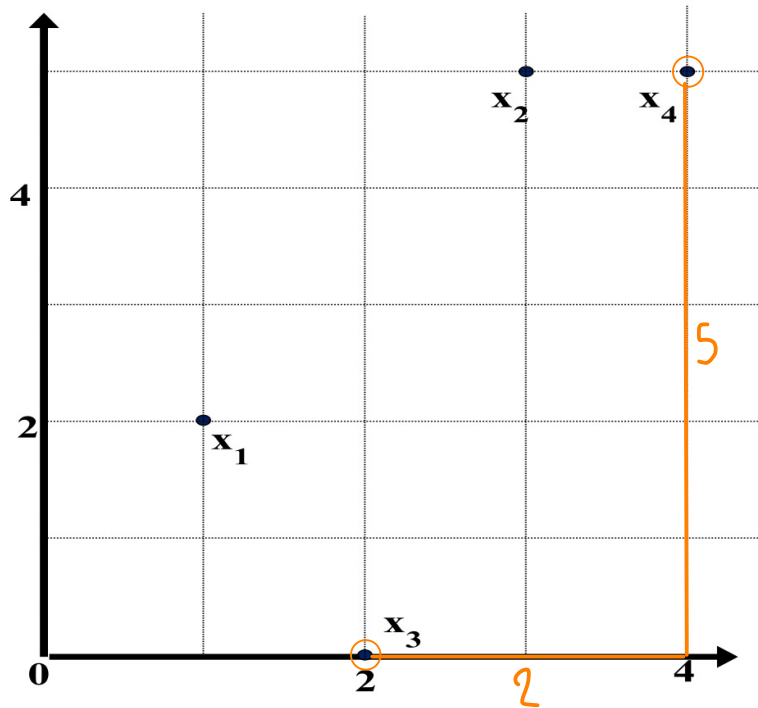
$\therefore c = \sqrt{13} = 3.61$

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

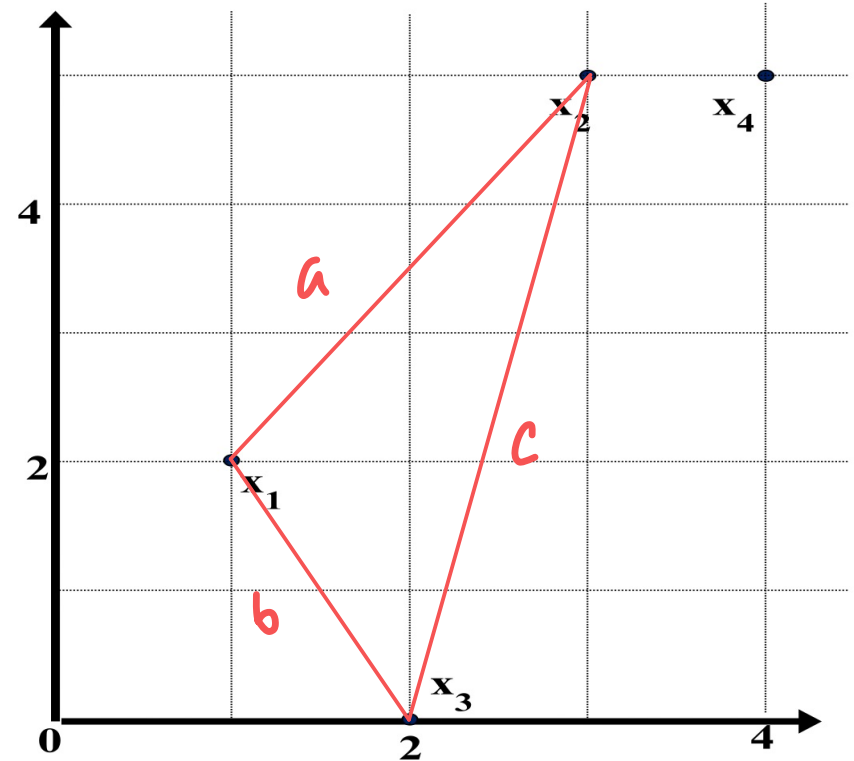


$L_1 \text{ norm}$ $= 3 + 3$ $= 6$	$L_2 \text{ norm}$ $= \sqrt{3^2 + 3^2}$ $= \sqrt{9 + 9}$ $= \sqrt{18}$	$L_\infty \text{ norm}$ $= 3$
--	---	----------------------------------

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



$L_1 \text{ norm}$ $= 2 + 5$ $= 7$	$L_2 \text{ norm}$ $= \sqrt{2^2 + 5^2}$ $= \sqrt{29}$	$L_\infty \text{ norm}$ $= 5$
--	---	----------------------------------



Linear Algebra

ระนาบ เป็น  $a + b$  ระนาบที่มีพื้นที่ ๒ มิติ