

ASU

Chapter 2

Measuring Data Similarity and Dissimilarity

- เป็นวิธีการ Data เพื่อเปรียบเทียบข้อมูลที่มีลักษณะคล้าย

1. Similarity measure or similarity function

- เป็นวิธีการหาความสัมพันธ์ที่มีลักษณะ / สอดคล้อง

- ครอบคลุมค่า $0, 1$; 0 = ไม่มีความคล้าย

1 = มีความคล้าย

2. Dissimilarity (or distance) measure

- เป็นวิธีการหาความแตกต่างระหว่างข้อมูล เช่น ระยะ 2 จุด บนเส้น

จำนวนจริงหรือจำนวนจริง $[0, 1]$ และ $[0, \infty]$

ใช้เปรียบเทียบข้อมูล

3. Proximity ความเป็นไปได้ที่ข้อมูลจะคล้ายกัน

การวัดระยะทาง ระยะทาง Distance matrix

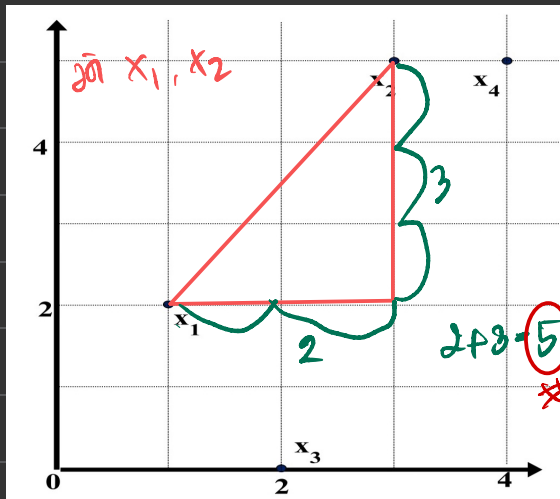
1. L_1 norm เป็นการวัดระยะแบบบล็อก

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

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

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

Ex //



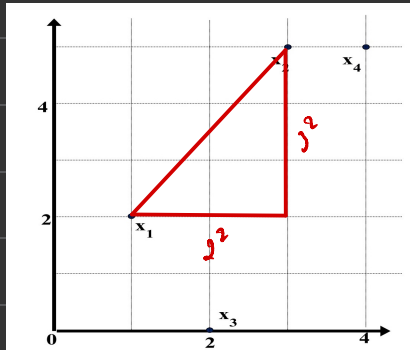
2. L_2 norm เป็นการวัดโดยวิธีสูตร มีทฤษฎี $A^2 + B^2 = C^2$

2. L_2 norm เป็นการวัดโดยใช้สูตร มีทริกซ์
 $A^2 + B^2, C^2$

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

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

E_x วัดจาก x_1, x_2



$$2^2 + 2^2 = 4 + 4$$

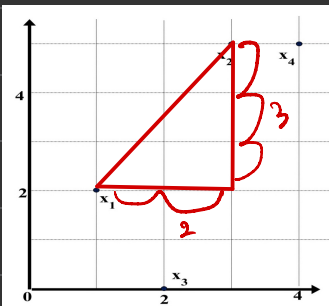
$$13 = 2^2 \cdot C = \sqrt{13}$$

3. L_∞ หรือ L_{\max} norm วัดขนาดที่จุด

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

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

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



เลือก จำนวนที่มากที่สุด คือ 2

Proximity for Binary

วิธีวัดความใกล้เคียง binary 2 วิธี

- Symmetric
- asymmetric

- A contingency table for binary data

		Object j		
		1	0	sum
Object 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:

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

วิธีวัด categorical 2 วิธี
1 simple matching

- Method 1: Simple matching

- m : # of matches, p : total # of variables

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

จำนวนตัวแปรที่ตรงกัน
หาร จำนวนตัวแปร

จำนวนตัวแปรทั้งหมด

2. วิธีวัดที่นิยมใน binary

Ordinal

ค่าของ z_{if} , 0, 1

ปัจจัย f

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

$i \in (1, \dots, M_f)$

M , ตัวชี้วัดค่าทั้งหมด.