

Similarity, Dissimilarity, and Proximity

- **Similarity measure** or similarity function \rightarrow ๑. ๒. ๓. ๔. ๕. ๖. ๗. ๘. ๙. ๑๐. \rightarrow ๐.๑.๒.๓.๔.๕.๖.๗.๘.๙.๑๐. $[0,1]$
 - 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 \Rightarrow ๑. ๒. ๓. ๔. ๕. ๖. ๗. ๘. ๙. ๑๐. \Rightarrow ๐.๑.๒.๓.๔.๕.๖.๗.๘.๙.๑๐.
 - 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) \rightarrow ๐
 - Range $[0, 1]$ or $[0, \infty)$, depending on the definition
- 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 \Rightarrow
- Dissimilarity (distance) matrix \rightarrow ๑. ๒. ๓. ๔. ๕. ๖. ๗. ๘. ๙. ๑๐.
 - 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

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

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

Standardizing Numeric Data การปรับ Scale ในมิติต่างๆ

- 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

<#>

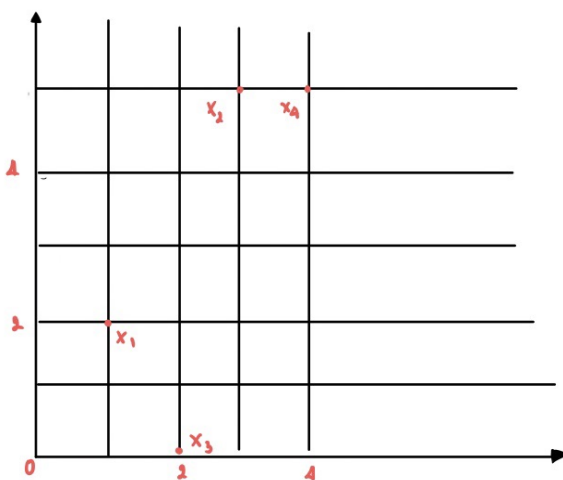
Example: Data Matrix and Dissimilarity Matrix

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 |



ค่าเดียวกัน

ค่าเดียวกัน

<#>

$\Rightarrow 3 \text{ MJV}$

- ๗ มนต์มหาเมฆ

[?]

- คัดจำผ่านของจุด

②

- อุณหภูมิมากขึ้น

1

 $\langle \# \rangle$