

“Measuring Data Similarity and Dissimilarity”

ဒိန္ဒကျော် ဦးသီရိ

b45020059-9



Similarity, Dissimilarity, and Proximity

- Similarity measure or similarity function

- Similarity ສືບຕະການເນື້ອຂອງ .Similarity Measure ທີ່ເປົ້າກວ່າຮອດການເນື້ອໂຄສົນໄວ້ 0-1 ເພື່ອໃຫ້ກວດສອບຜົນດໍາລັດຂອງ

ถ้าเน้นเรื่องความต่อเนื่องใน \mathbb{R}^n ดูว่า function ในที่นี่ก็ต้องต่อเนื่องทุกจุด แต่จะต้องมีค่าอนันต์

- ## • Dissimilarity (or distance)

- การวัดว่าตัวอย่างคุณภาพใดในฝาดีเจ ต้องอยู่ distance 有多遠 ใจความของคุณภาพนั้นๆ ใจความ similarity ใจความ inverse ใจความ

- Proximity

- Proximity តើអាជីវកម្មនេះទាំងនេះទូទៅដោយខ្លួន

Data Matrix and Dissimilarity Matrix

Data Matrix

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

- ວິທານໍາແກ້ໄຂ ດ້ວຍ Matrix ກີ່ສື່ distance matrix

$$\left(\begin{array}{cccc} 0 & & & \\ d(2,1) & 0 & & \\ \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & 0 \end{array} \right)$$

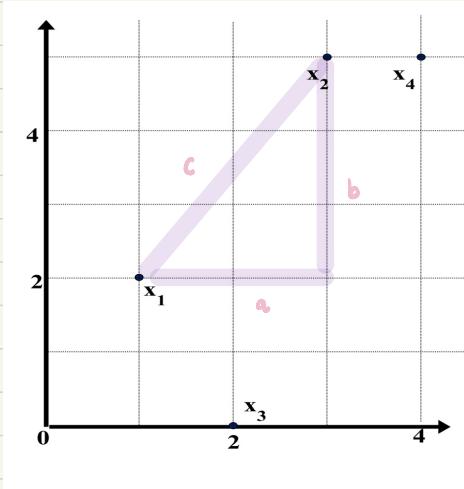
နိတ်ဆက်ချေ အ ရွှေ မီတာ မီတာ မီတာ မီတာ

Standardizing Numeric Data

ពំលេចទិន្នន័យការងារសាធារណរដ្ឋប្រជាជាតិ និងក្រសួងពេទ្យ នៃរដ្ឋបាល នៅថ្ងៃទី ២០ ខែ មីនា ឆ្នាំ ២០១៩ (០,១)

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

Example: Data Matrix and Dissimilarity



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

ទីផ្សារក្នុងការគិតវិភាគ $C = a^2 + b^2$

ទីផ្សារក្នុងការគិតវិភាគ x_1 និង x_2 ពេលវិភាគ $C = 3^2 + 3^2 \Rightarrow C = 18$

Minkowski distance

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

រូបំពេល 1 (L_1 norm) Manhattan (or city block) distance

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

(ទីផ្សារក្នុងការគិតវិភាគដែលមិនស្ថិតនៅលើលេខគិតគុណ)

រូបំពេល 2 (L_2 norm) Euclidean distance

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

(ទីផ្សារក្នុងការគិតវិភាគដែលមិនស្ថិតនៅលើលេខគិតគុណ)

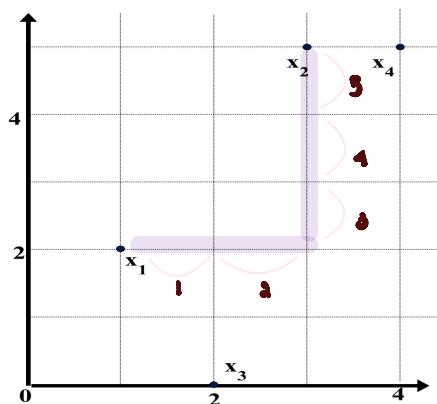
$p = \infty$: (L_∞ norm, L_∞ norm) "supremum" distance

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

Example: Minkowski Distance at Special Cases

L₁

point	attribute 1	attribute 2
x ₁	1	2
x ₂	3	5
x ₃	2	0
x ₄	4	5



Manhattan (L₁)

L	x ₁	x ₂	x ₃	x ₄
x ₁	0			
x ₂	5	0		
x ₃	3	6	0	
x ₄	6	1	7	0

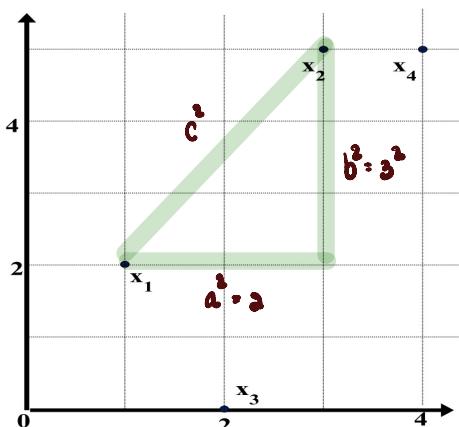
ចំណាំរាយនៃទីលក្ខណន៍ នៅលើ មាលបុង ហើយ x₁ តិច y₂ ដោយពេល 5 ទី

វិវឌ្ឍន៍ x₁ ដើម្បី 2 នូវ និង 3 នូវ និង 4 នូវ

ចំណាំ x₁ និង x₂ ដោយពេល 2+3 = 5 ទី

L₂

point	attribute 1	attribute 2
x ₁	1	2
x ₂	3	5
x ₃	2	0
x ₄	4	5



Euclidean (L₂)

L ₂	x ₁	x ₂	x ₃	x ₄
x ₁	0			
x ₂	3.61	0		
x ₃	2.24	5.1	0	
x ₄	4.24	1	5.39	0

កិច្ចការស្ថិកការណ៍ c = a² + b²

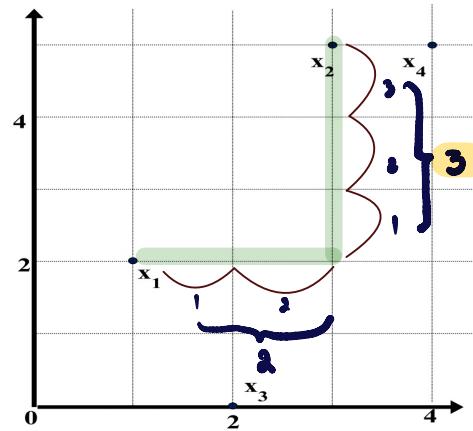
$$c = 2^2 + 3^2$$

$$c = \sqrt{1+9}$$

ចំណាំ C = 3.61

L_∞

point	attribute 1	attribute 2
x ₁	1	2
x ₂	3	5
x ₃	2	0
x ₄	4	5



Supremum (L_∞)

L _∞	x ₁	x ₂	x ₃	x ₄
x ₁	0			
x ₂	3	0		
x ₃	2	5	0	
x ₄	3	1	5	0

ຈະກຳຕ່າງໆ Max ຍົວວ່າຍຸດ ເນັ້ນພາຍາການນໍາກວ່າ
x₁ ດັ່ງ x₂ ໂດຍ x₁ ດີວຽວນໍາກວ່າ x₂ ແລະ ທາງຫຼຸດ x₂ ເປັນ
ກັບ x₃ ອູດ
ຢູ່ວ່າຜົນ L_∞ ສົ່ງໄວ ເນັ້ນຕະຫຼາມ Max

Proximity Measure for Binary Attributes

Symmetric binary (ສິ້ນມະນີ້ນຕື່ມເຕັກ)

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

		Object j	
		1	0
Object i	1	r	s
	0	t	s+t
sum		r+s	r+t

Asymmetric binary (ຄອບຄ່າທີ່ມີການແຕ່ງຕ່າງ)

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

Example: Dissimilarity between Asymmetric Binary Variables

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

Tayari 4 wa. P=2 wa. N=0

		Mary		
		1	0	Σ_{row}
Jack		1	2	0
0	1	1	3	4
	Σ_{col}	3	3	6
lim				

ເນື້ອມສົງໄນ້ກຳນົດ ແລະ ອົບເວັບ ຢ. ສູນເວັບ ແລະ ຢ. ສູນເວັບ

កំណត់សាលាបុរីជាមួយនឹងការបិទសាលាបុរី (asymmetric binary)

$$\text{Ans: } d(\text{jcek, Mary}) = \frac{r+s}{q+r+s}$$
$$= \frac{0+1}{2+0+1}$$
$$= 0.33$$

Proximity Measure for Categorical Attributes

$$d(i,j) = \frac{P-m}{P} ; P: total \rightarrow 3 (ສັນໄງ້ ຂອງລົດ ຄຸນວັດ ສັກຄາ ອົງກ)$$

m: matches

	កុំភាគ	សេចក្តី	លក្ខណៈ
កំ	កុំភាគ	កុំភាគ	កុំភាគ
កំ	កុំភាគអ៊ូនិស្ស	កុំភាគ	Programmer
កំ	កុំភាគ	កុំភាគអ៊ូនិស្ស	Programmer
កំ	កុំភាគអ៊ូនិស្ស	កុំភាគ	កុំភាគអ៊ូនិស្ស

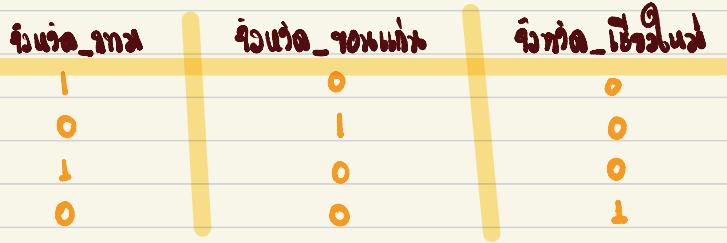
distances matrix = 4×4

	d_1	d_2	d_3	d_4
d_1	0			
d_2	1	0		
d_3	0.67		0	
d_4				0

$$\therefore d(d_1, d_2) = \frac{3 \cdot 0}{2}$$

$$\therefore d(m_1, m_2) = \frac{3-1}{3} = \frac{2}{3} = 0.667$$

dummy



Ordinal Variables

เป็น เวิร์ชั่นของค่านักเรียน

$$z_{ij} = \frac{r_{ij}-1}{n_j-1} \quad ; r_{ij} \in \{1, 2, \dots, n\}$$

Example $r_{ij} \in \{1, \dots, 4\}$

freshman:1 Sophomore:2 junior:3 senior:4

$$\begin{array}{cccc} z = 1 - 1 & z = 2 - 1 & z = 3 - 1 & z = 4 - 1 \\ 4-1 & 4-1 & 4-1 & 4-1 \end{array}$$

ค่าต่อไปนี้คือค่าที่ได้มาจากการคำนวณ

$$d(junior, senior) = 1 - \frac{3}{3} + \frac{1}{3}$$

