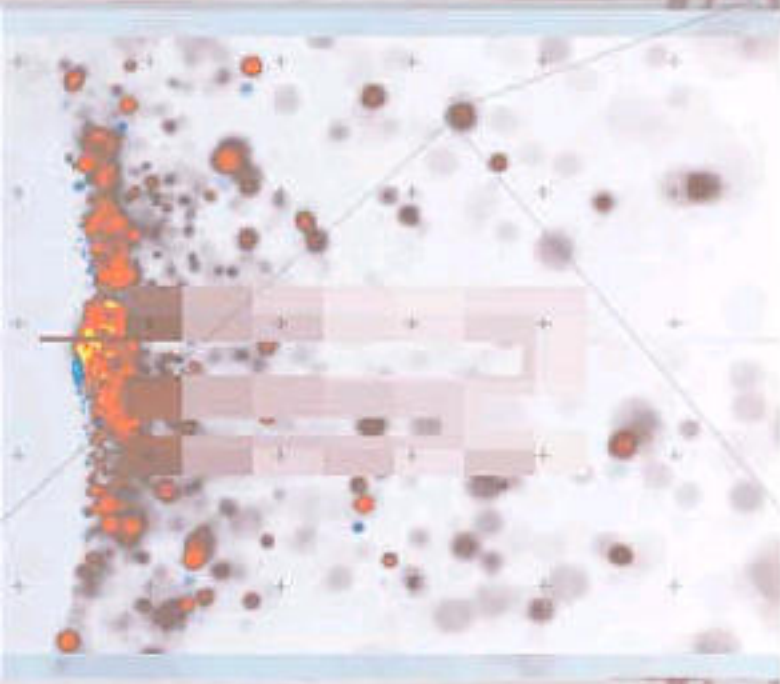


The background features a complex geometric pattern of thin, intersecting lines in shades of brown and grey, creating a mesh-like effect. Scattered throughout are small, colored dots in green, blue, and orange. A semi-transparent white banner with a subtle geometric shape is positioned across the upper middle of the image.

Distance on Numeric Data: Minkowski Distance



Data Matrix and Dissimilarity Matrix

□ Data matrix

- A data matrix of n data points with l dimensions

$n \times l$
数据量 维度
→

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

□ Dissimilarity (distance) matrix

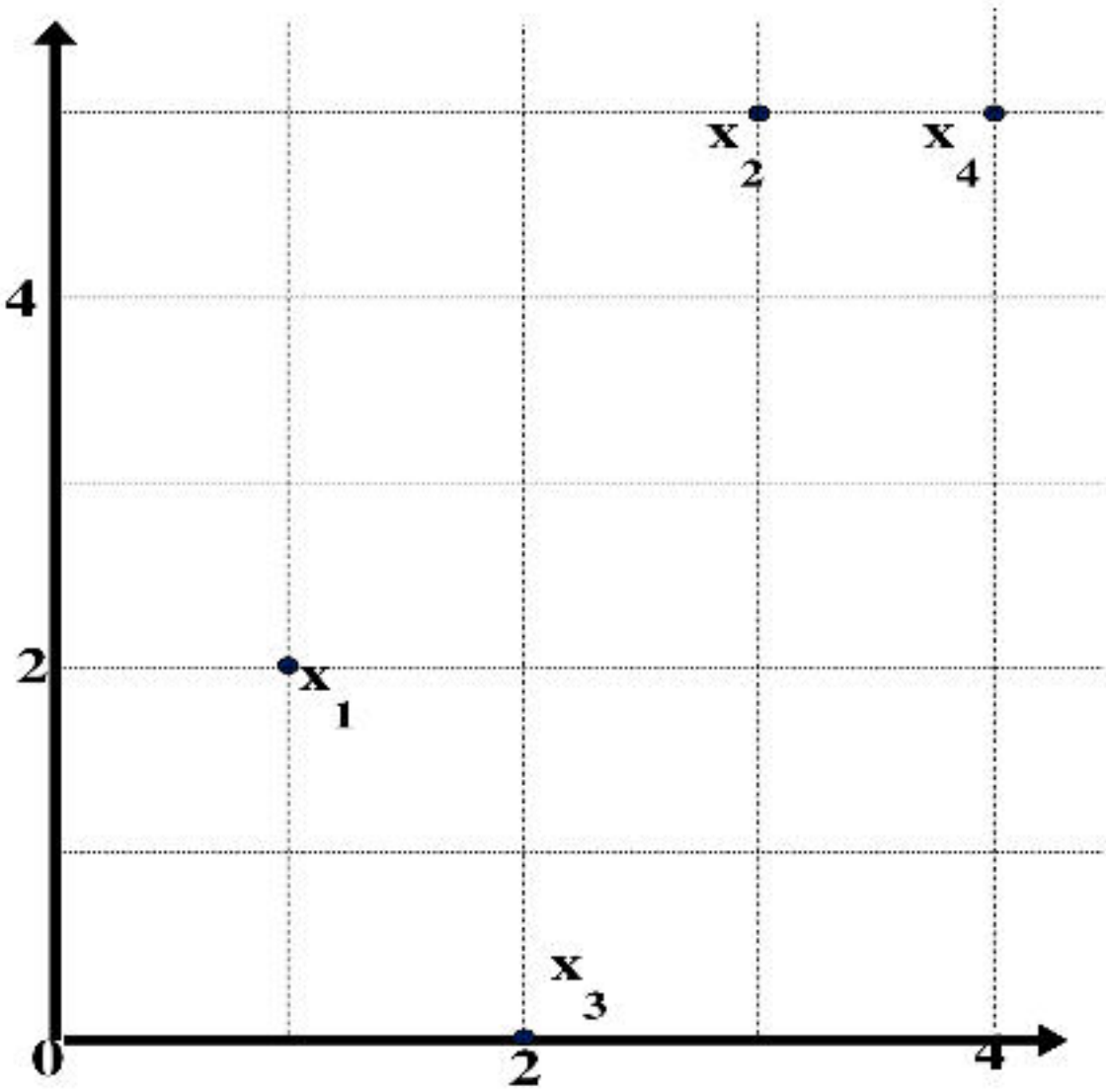
- n data points, but registers only the distance $d(i, j)$ (typically metric)
- Usually symmetric, thus a triangular matrix
- **Distance functions** are usually different for real, boolean, categorical, ordinal, ratio, and vector variables
- Weights can be associated with different variables based on applications and data semantics



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

上三角阵

Example: Data Matrix and Dissimilarity Matrix



Data Matrix

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

$$d(x_i, x_j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2}$$

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**

- ☐ $d(i, j) > 0$ 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

三角形两边之和 > 第三边.



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

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

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

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

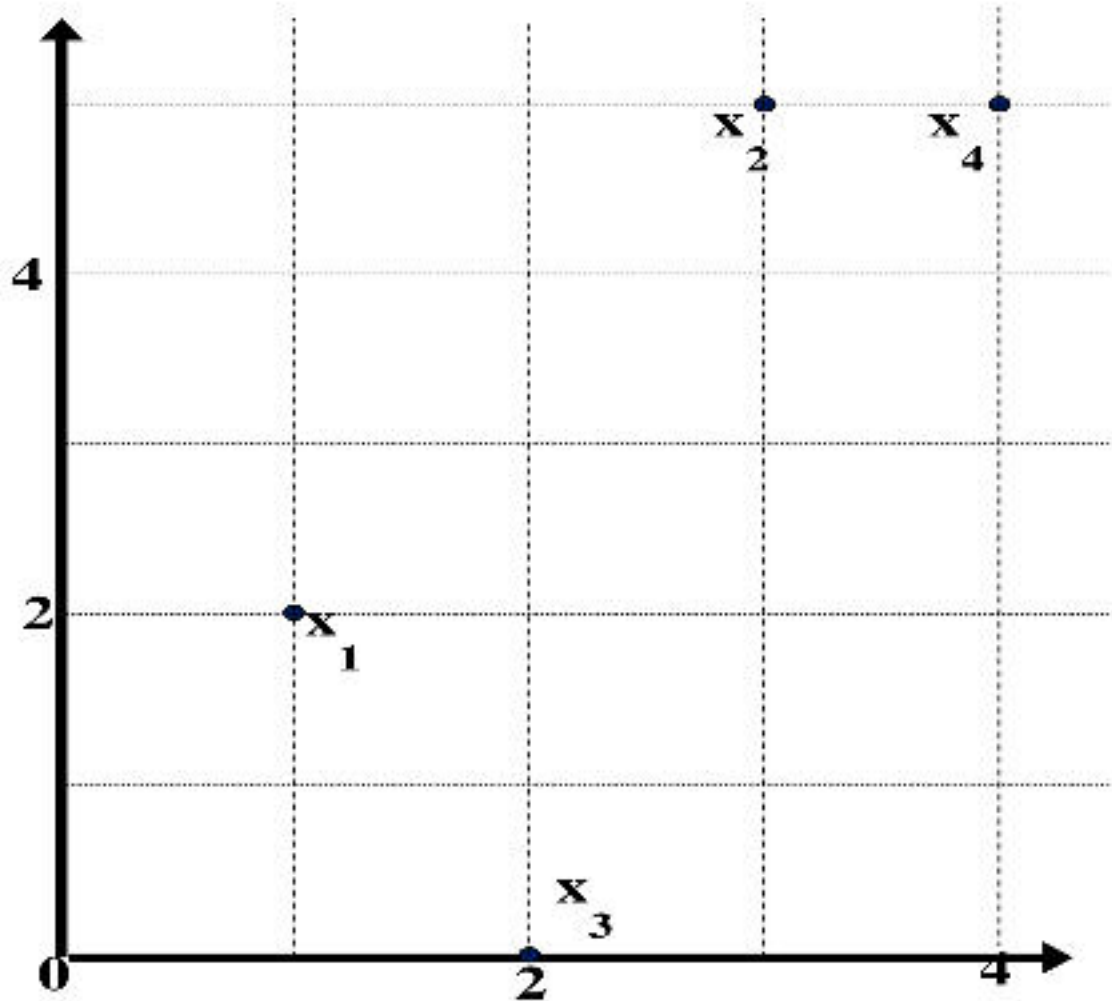
□ $p \rightarrow \infty$: (L_{\max} norm, L_{∞} norm) **"supremum" distance**

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

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

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_∞) \rightarrow 看 Manhattan 取最大的 Manhattan 距离.

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

(x1, x2)
|y|=3, |x|=2.
取3