

# Week 3

# Engineering Data (Part II)

Seokho Chi

Assistant Professor | Ph.D.

SNU Construction Innovation Lab



Source: Tan, Kumar, Steinback (2006)

# Similarity and Dissimilarity

- Similarity
  - Numerical measure of how alike two data objects are.
  - Is higher when objects are more alike.
  - Often falls in the range  $[0,1]$
- Dissimilarity
  - Numerical measure of how different are two data objects
  - Lower when objects are more alike
  - Minimum dissimilarity is often 0
  - Upper limit varies
  - **Distance**: special class of dissimilarity
- Proximity refers to a similarity or dissimilarity

# Similarity/Dissimilarity for Simple Attributes

$p$  and  $q$  are the attribute values for two data objects.

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$	$s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$
Ordinal	$d = \frac{ p-q }{n-1}$ (values mapped to integers 0 to $n-1$ , where $n$ is the number of values)	$s = 1 - \frac{ p-q }{n-1}$
Interval or Ratio	$d =  p - q $	$s = -d, s = \frac{1}{1+d} \text{ or } s = 1 - \frac{d - \min\_d}{\max\_d - \min\_d}$

**Table 5.1.** Similarity and dissimilarity for simple attributes

# Euclidean Distance

(Dissimilarity)

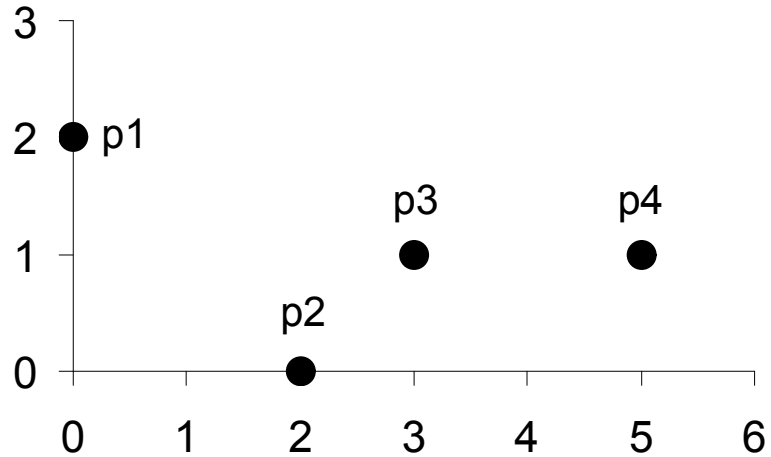
- Euclidean Distance (distance b/w points)

$$\mathit{dist} = \sqrt{\sum_{k=1}^n (p_k - q_k)^2}$$

Where  $n$  is the number of dimensions (attributes) and  $p_k$  and  $q_k$  are, respectively, the  $k^{\text{th}}$  attributes (components) or data objects  $p$  and  $q$ .

- Standardization is necessary, if scales differ.

# Euclidean Distance



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

**Distance Matrix**

# Minkowski Distance

- Minkowski Distance is a generalization of Euclidean Distance

$$\mathit{dist} = \left( \sum_{k=1}^n |p_k - q_k|^r \right)^{\frac{1}{r}}$$

Where  $r$  is a parameter,  $n$  is the number of dimensions (attributes) and  $p_k$  and  $q_k$  are, respectively, the  $k$ th attributes (components) or data objects  $p$  and  $q$ .

# Minkowski Distance: Examples

- $r = 1$ . City block (Manhattan, taxicab,  $L_1$  norm) distance.
  - A common example of this is the Hamming distance, which is just **the number of bits** that are different between two binary vectors
- $r = 2$ . Euclidean distance
- $r \rightarrow \infty$ . “supremum” ( $L_{\max}$  norm,  $L_{\infty}$  norm) distance.
  - This is the maximum difference between **any component of the vectors**
- Do not confuse  $r$  with  $n$ , i.e., **all these distances are defined for all numbers of dimensions “n”.**

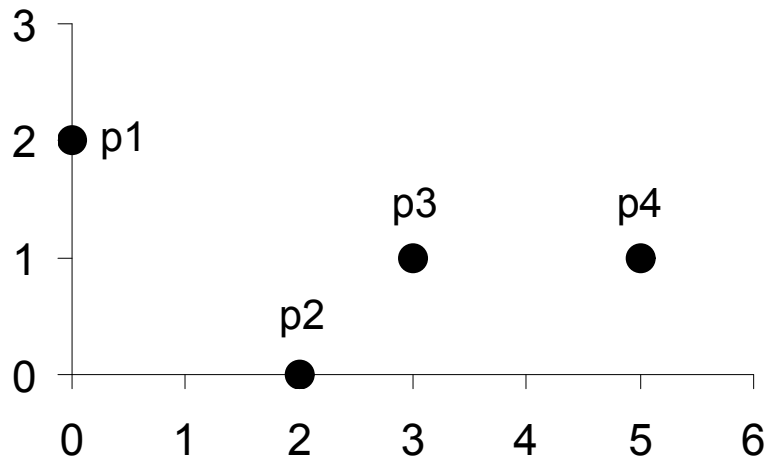
# Minkowski Distance

point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

L1	p1	p2	p3	p4
p1	0	4	4	6
p2	4	0	2	4
p3	4	2	0	2
p4	6	4	2	0

L2	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

L <sub>∞</sub>	p1	p2	p3	p4
p1	0	2	3	5
p2	2	0	1	3
p3	3	1	0	2
p4	5	3	2	0

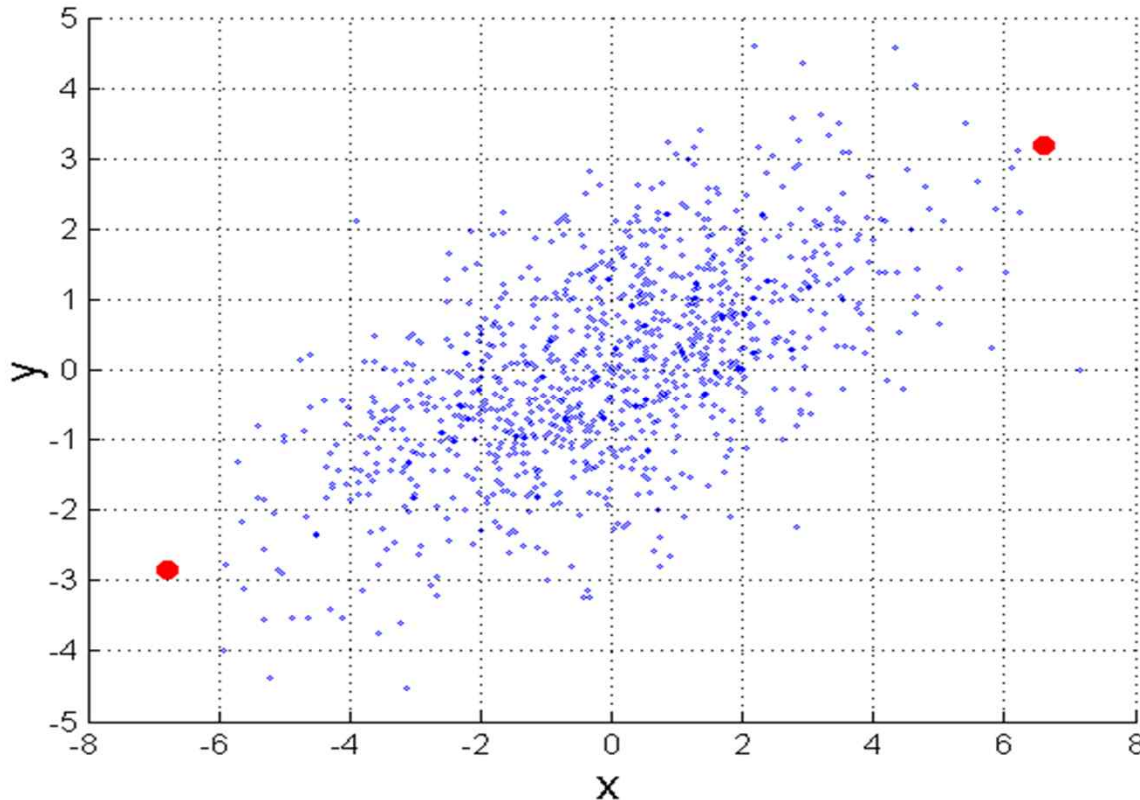


**Distance Matrices**



# Mahalanobis Distance

$$\mathit{mahalanobis}(p, q) = (p - q) \Sigma^{-1} (p - q)^T$$



**Distance b/w the point and the distribution mean**

***X times error than SD***

*(평균과의 거리가 표준편차의 몇 배인가)*

**$\Sigma$  is the covariance matrix of the input data  $X$**

$$\Sigma_{j,k} = \frac{1}{n-1} \sum_{i=1}^n (X_{ij} - \bar{X}_j)(X_{ik} - \bar{X}_k)$$

**For red points, the Euclidean distance is 14.7, Mahalanobis distance is 6.**

*Explain how difficult it occurs or how strange the point is: Outlier detection*

*교통량 20대에 표준편차 3대일 경우,*

*26대가 지나가면 평균과의 거리는 6이지만 Mahalanobis distance는 6/3=2 즉, 표준적인 편차의 2배정도의 오차*

# Common Properties of a Distance

- Distances, such as the Euclidean distance, have some well known properties.

1.  $d(p, q) \geq 0$  for all  $p$  and  $q$

$d(p, q) = 0$  only if  $p = q$  (Positive definiteness)

2.  $d(p, q) = d(q, p)$  for all  $p$  and  $q$  (Symmetry)

3.  $d(p, r) \leq d(p, q) + d(q, r)$  for all points  $p, q,$  and  $r$   
(Triangle Inequality)

where  $d(p, q)$  is the distance (dissimilarity) between points (data objects)  $p$  and  $q$

# Common Properties of a Similarity

- Similarities, also have some well known properties.

1.  $s(p, q) = 1$  (or maximum similarity) only if  $p = q$

2.  $s(p, q) = s(q, p)$  for all  $p$  and  $q$  (Symmetry)

where  $s(p, q)$  is the similarity between points (data objects)  $p$  and  $q$

# Similarity Between Binary Vectors



- Common situation is that objects,  $p$  and  $q$ , have only binary attributes

- Compute similarities using the following quantities

$M_{01}$  = the number of attributes where  $p$  was 0 and  $q$  was 1

$M_{10}$  = the number of attributes where  $p$  was 1 and  $q$  was 0

$M_{00}$  = the number of attributes where  $p$  was 0 and  $q$  was 0

$M_{11}$  = the number of attributes where  $p$  was 1 and  $q$  was 1

- Simple Matching and Jaccard Coefficients

SMC = number of matches / number of attributes

$$= (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00})$$

J = number of 11 matches / number of not-both-zero attributes values

$$= (M_{11}) / (M_{01} + M_{10} + M_{11})$$

→ *Ignore 0-0 matches to avoid miss-matches by noisy 0 values*

# Cosine Similarity

- If  $d_1$  and  $d_2$  are two document vectors, then

$$\cos(d_1, d_2) = (d_1 \bullet d_2) / \|d_1\| \|d_2\| ,$$

where  $\bullet$  indicates vector dot product and  $\|d\|$  is the length of vector  $d$ .

*Jaccard measure + non-binary vectors*

- Example:

$$d_1 = \mathbf{3\ 2\ 0\ 5\ 0\ 0\ 0\ 2\ 0\ 0}$$

$$d_2 = \mathbf{1\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 2}$$

$$d_1 \bullet d_2 = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$$

$$\|d_1\| = (3*3 + 2*2 + 0*0 + 5*5 + 0*0 + 0*0 + 0*0 + 2*2 + 0*0 + 0*0)^{0.5} = (42)^{0.5} = 6.481$$

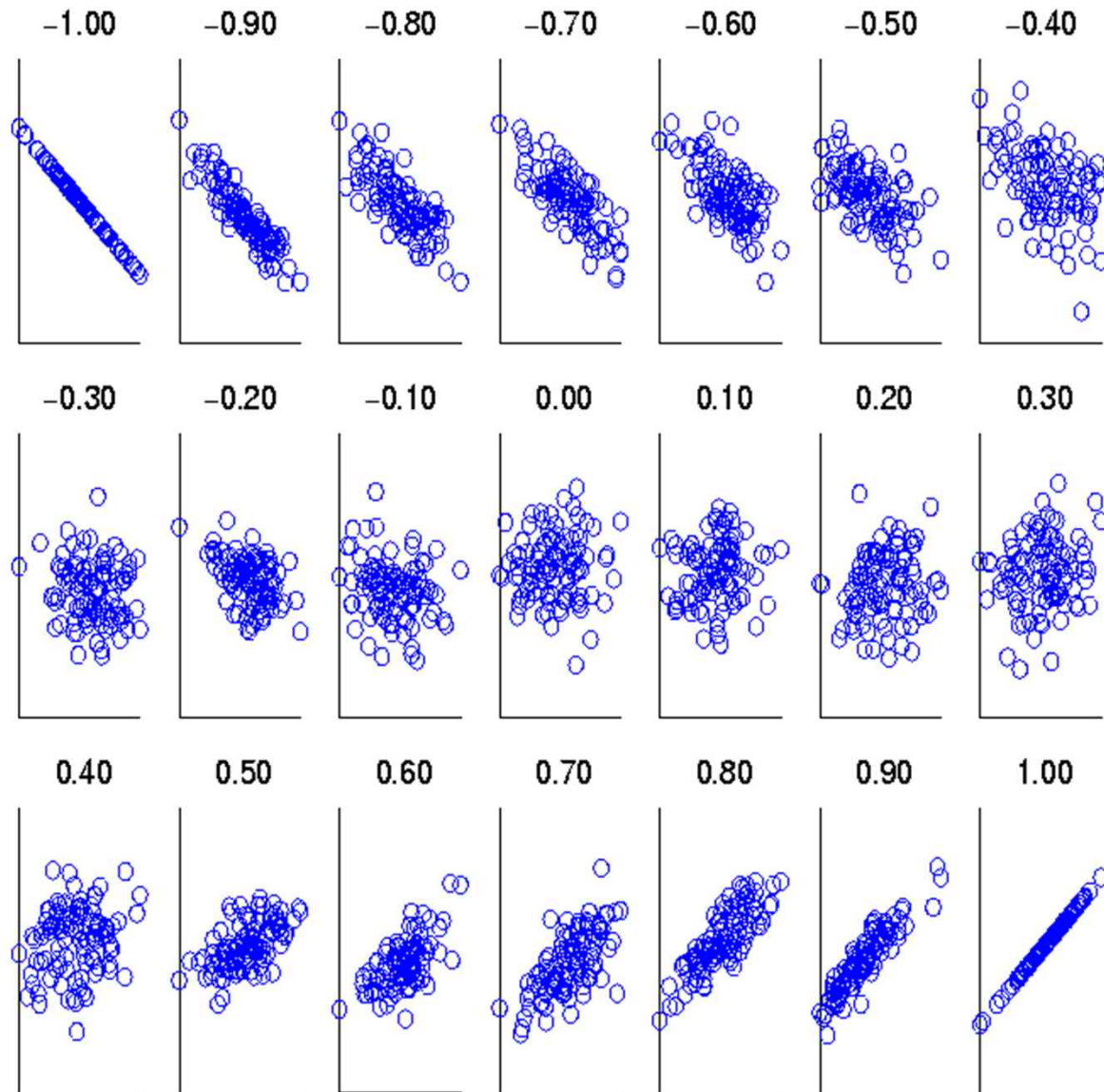
$$\|d_2\| = (1*1 + 0*0 + 0*0 + 0*0 + 0*0 + 0*0 + 0*0 + 1*1 + 0*0 + 2*2)^{0.5} = (6)^{0.5} = 2.245$$

$$\cos(d_1, d_2) = .3150 \quad (1 \rightarrow 0^\circ \rightarrow \text{same except length}; 0 \rightarrow 90^\circ \rightarrow \text{do not share})$$

# Correlation

- Correlation measures the linear relationship between objects
- To compute correlation, we standardize data objects,  $p$  and  $q$ , and then take their dot product

# Visually Evaluating Correlation



**Scatter plots showing the similarity from -1 to 1.**

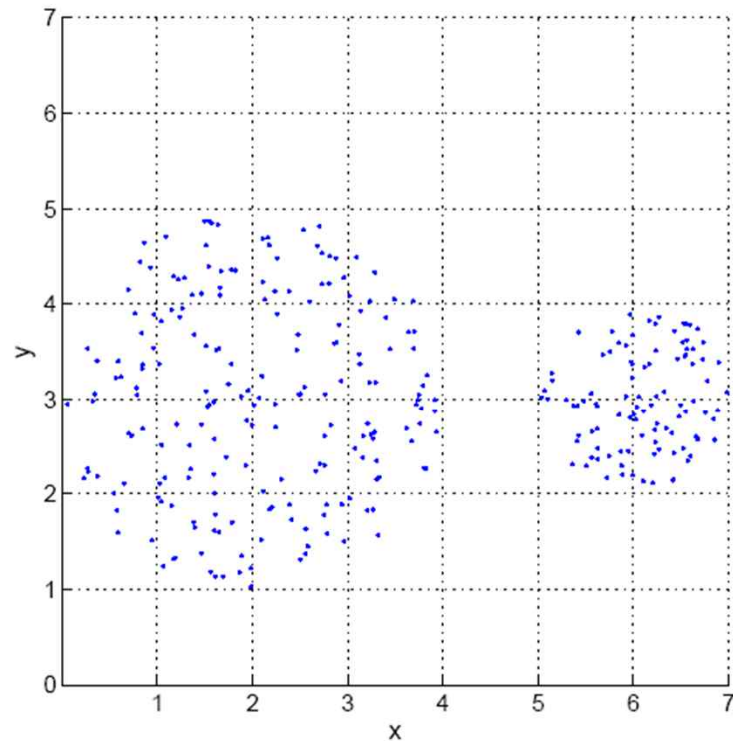
# Density

- Density-based clustering require a notion of density
- Examples:
  - Euclidean density
    - Euclidean density = number of points per unit volume
  - Probability density
  - Center-based density



# Euclidean Density – Cell-based

- Simplest approach is to divide region into a number of rectangular cells of equal volume and define density as # of points the cell contains



**Figure 7.13.** Cell-based density.

0	0	0	0	0	0	0
0	0	0	0	0	0	0
4	17	18	6	0	0	0
14	14	13	13	0	18	27
11	18	10	21	0	24	31
3	20	14	4	0	0	0
0	0	0	0	0	0	0

**Table 7.6.** Point counts for each grid cell.

# Euclidean Density – Center-based

- Euclidean density is the number of points within a specified radius of the point

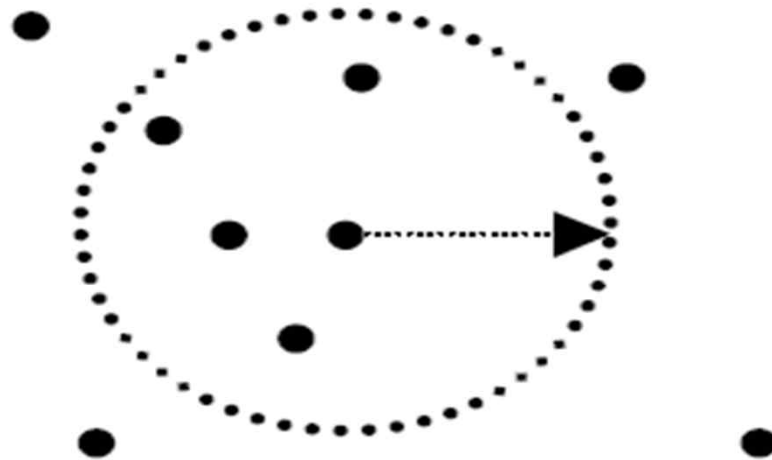


Figure 7.14. Illustration of center-based density.