

Chapter 4 Clustering

Associate Professor Yachai Limpiyakorn, Ph.D.
Department of Computer Engineering
Chulalongkorn University

Natural Grouping

- Unsupervised learning

- Natural grouping:

- High intra-class similarity
- Low inter-class similarity

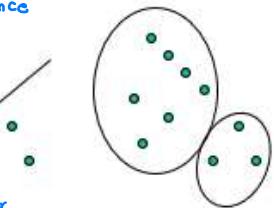
- As a tool for understanding data distribution or data preprocessing prior to DM modeling

- .preprocessing outlier

- Clustering is subjective

ເນື້ອຍຸປົກການຝຶກ, ຕົວຢັດ

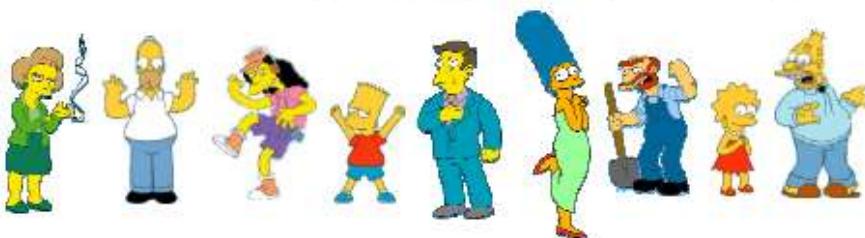
ເທົ່ານັ້ນ, ຈີງກຸ່ມຂາວາງລົກປະບຸ



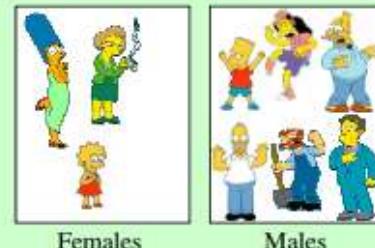
2110773-4 2/66

2

What is a natural grouping among these objects?



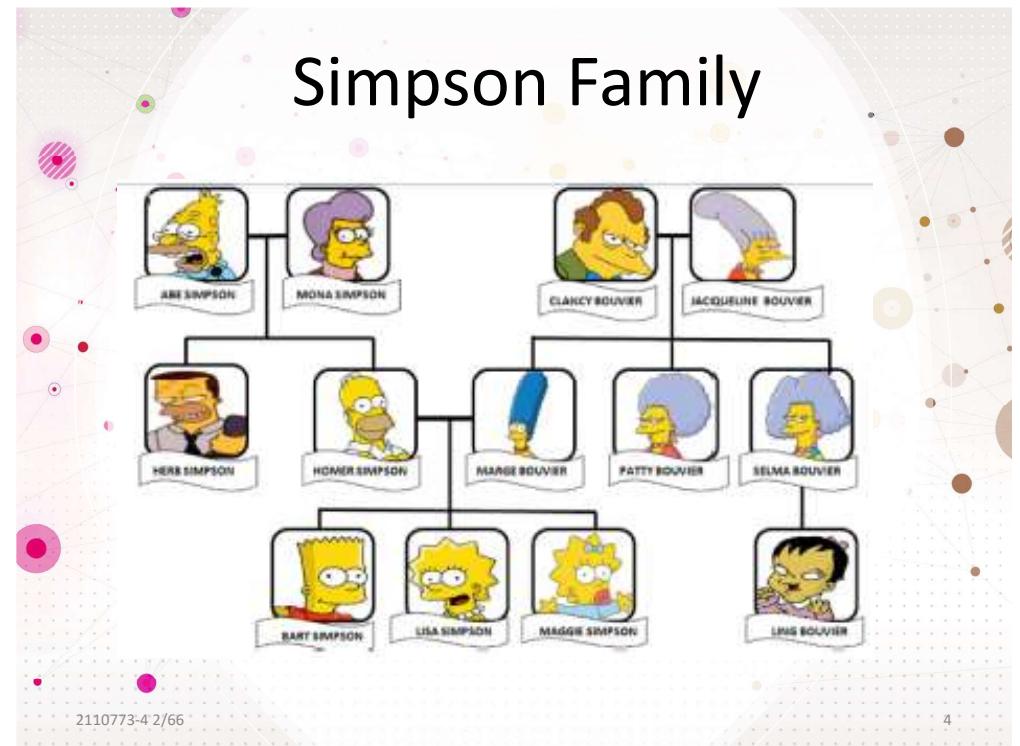
Clustering is subjective



2110773-4 2/66

3

Simpson Family



2110773-4 2/66

4

What is similarity?

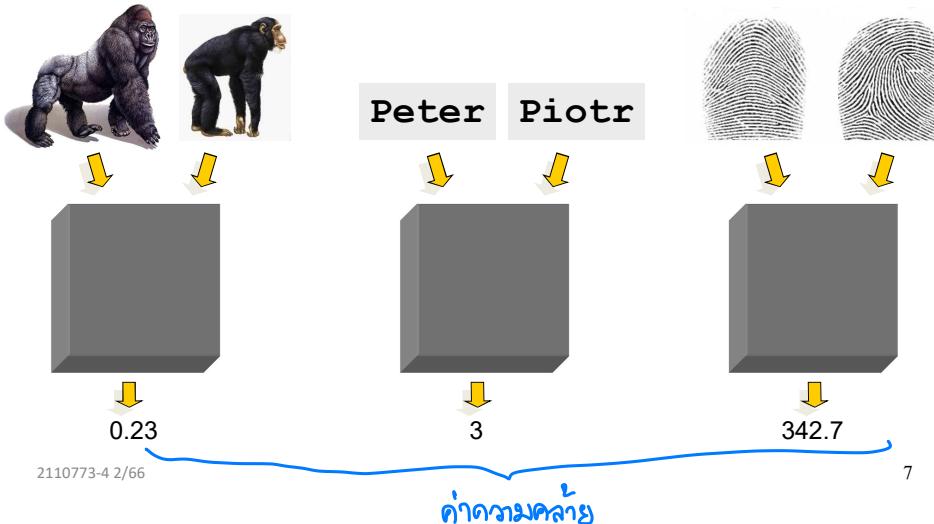


2110773-4 2/66

5

Defining distance measures

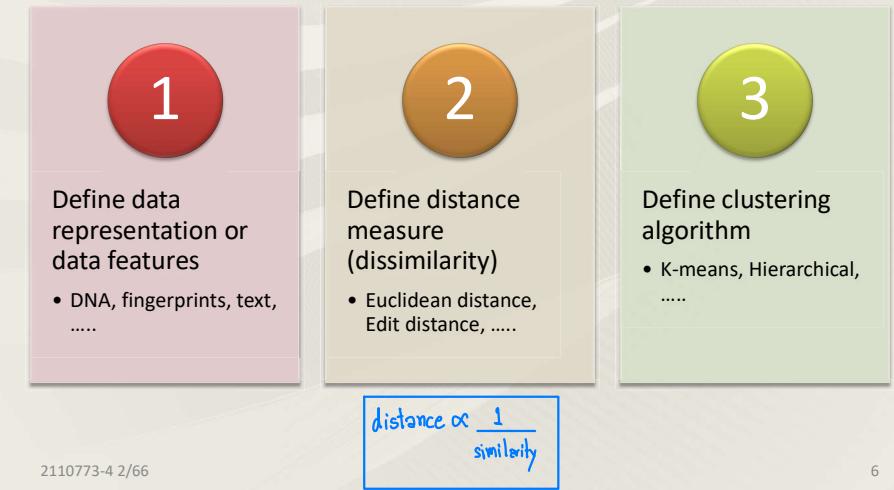
Definition: Let O_1 and O_2 be two objects from the universe of possible objects. The distance (dissimilarity) between O_1 and O_2 is a real number denoted by $D(O_1, O_2)$



2110773-4 2/66

7

Before Clustering



2110773-4 2/66

6

Edit Distance Example

It is possible to transform any string Q into string C , using only *Substitution*, *Insertion* and *Deletion*. Assume that each of these operators has a cost associated with it.

The similarity between two strings can be defined as the cost of the cheapest transformation from Q to C .

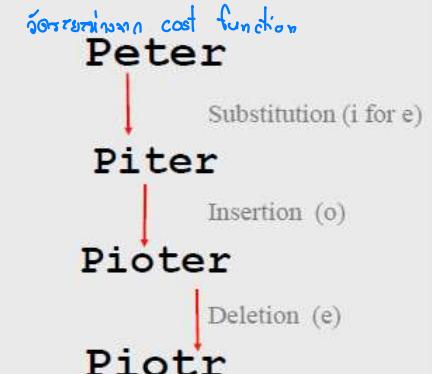
Note that for now we have ignored the issue of how we can find this cheapest transformation.

How similar are the names "Peter" and "Piotr"?

Assume the following cost function

Substitution	1 Unit
Insertion	1 Unit
Deletion	1 Unit

$D(\text{Peter}, \text{Piotr})$ is 3



2110773-4 2/66

8

Distance Measures

- Edit/ Transformation distance
 - $D(\text{Patty}, \text{Selma}) = 4$
 - Change dress color 1 pt.
 - Change shoe color 1 pt. มากที่สุดที่ต้องเปลี่ยน 4 อย่าง
 - Change earring shape 1 pt.
 - Change hair part 1 pt.
- Minkowski distance
 - Manhattan
 - Euclidean



2110773-4 2/66

9

Properties of Distance Measure

$$D(A, B) \geq 0$$

Positivity

$$D(A, B) = D(B, A)$$

Symmetry

$$D(A, B) = 0 \text{ iff } A = B$$

Reflexive

$$D(A, B) \leq D(A, C) + D(B, C)$$

Triangle Inequality

2110773-4 2/66

11

Minkowski Distance

- Some popular ones include: *Minkowski distance*:

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

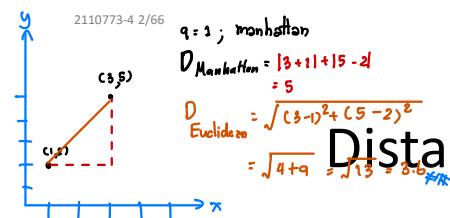
where $i = (x_{i1}, x_{i2}, \dots, x_{ip})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jp})$ are two p -dimensional data objects, and q is a positive integer

- Manhattan distance

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

- Euclidean distance ผลลัพธ์ตามมีปรับแต่ง

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



Distance Matrix

ข้อมูลจะถูกจัดเก็บในเมตริกซ์ (two mode)

ประกอบด้วยข้อมูล ก ตัว แต่ละตัวมี

คุณลักษณะ p คุณลักษณะ

f บน column attribute

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{np} \end{bmatrix}_{n \times p}$$

f ตัว 1 ถึง p

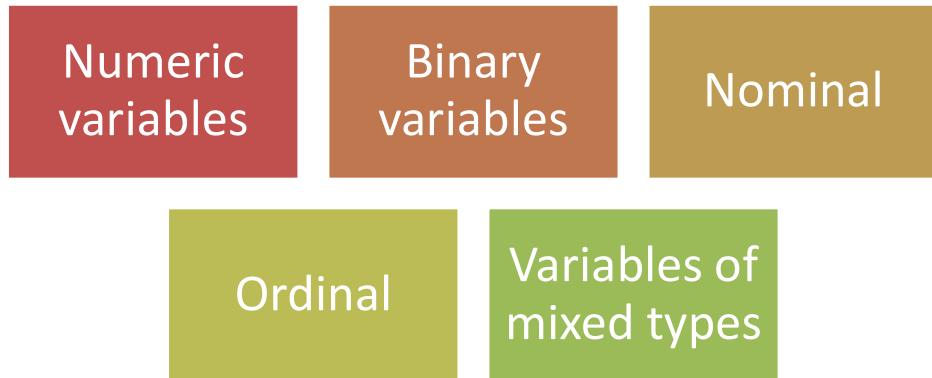
มาตรฐานความ(ไม่)คล้ายอื่นในรูปแบบ
ฟังก์ชันระยะห่าง $d(i, j)$ ซึ่งจัดเก็บในเมตริกซ์
ความไม่คล้าย นิยามฟังก์ชันระยะห่าง $d(i, j)$
มักแต่ต่างกันไปตามประเภทข้อมูล

$$\begin{bmatrix} 0 & d(1,2) & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

10

12

ประเภทตัวแปร



2110773-4 2/66

13

Numeric Variable



Normalization or Standardization is suggested to transform to equivalent ranges
ทำให้ตัวแปรมีรูปแบบเดียวกัน



ตัวแปรที่ growth rate ไม่เป็นเชิงเส้น เช่น การเริ่มต้นโดยช่องแบคทีเรีย หรือ การสื่อสารด้วยของลูกศิริน์มีนักภาษาเพียงสี่คน take log() prior to Normalization

2110773-4 2/66

14

Binary Variable

Symmetric

- Simple Matching Coefficient

$$d(i,j) = \frac{b+c}{a+b+c+d}$$

หาตัวอย่าง

Asymmetric

- Jaccard Coefficient

$$d(i,j) = \frac{b+c}{a+b+c}$$

หาตัวอย่าง
ก็จะได้ 0 ถ้า 0 ไม่ใช่ค่า 1 ถ้า 1

c คือ attribute (0, 1)

obj			Object j			
	1	0		1	0	sum
i	f_1	f_2		a	b	$a+b$
j	0	1	1	c	d	$c+d$
k	1	0	0	$a+c$	$b+d$	p
l	1	0	0	$d(j,h)$	0	หน่วยค่า attribute
m	Th	Th		$d(i,j) = 4$	ต่างกัน 4 ค่า attribute	
n	BK	BK				
on nominal		non ordinal				

15

Nominal Variable

- วิธีการคำนวณความ(ไม่)คล้ายของตัวแปรที่รู้ค่านับ อาจเลือกใช้ฟังก์ชันระยะห่าง

- วิธีที่1: Simple Matching
กำหนดให้ m เป็นจำนวนตัวแปรที่รู้ค่านับที่วัดถูกทั้งคู่มีค่าเหมือนกัน

p เป็นจำนวนตัวแปรทั้งหมด ระยะห่างระหว่างวัดถูกทั้งสองค่าได้โดย

- วิธีที่2: แทนตัวแปรที่รู้ค่านับที่มีค่าที่เป็นไปได้ M ค่า ด้วยตัวแปรที่จำนวน M ตัว

ใช้ 7 bit กรณี 7 ค่า
↓
one hot

2110773-4 2/66

16

order 1, 2, 3
medal = {G, S, B}

Ordinal Variable

order

- การคำนวณหาความ(ไม่)คล้ายของวัตถุที่มีตัวแปรแบบมีอันดับ สามารถใช้พังก์ชันระยะห่างที่ใช้สำหรับตัวแปร interval-scaled ได้ โดยการแทนค่าตัวแปร interval-scaled ด้วยค่าอันดับของตัวแปร x_{if}

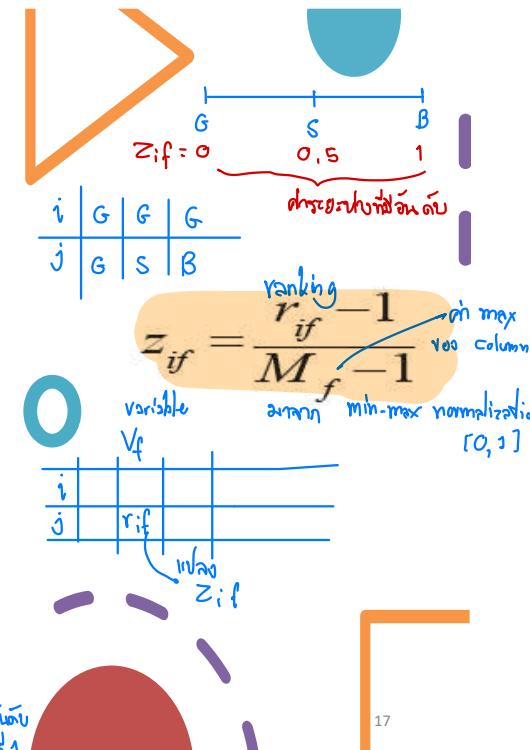
$$r_{if} \in \{1, \dots, M_f\}$$

- แปลงค่าตัวแปรมีอันดับที่ j^{th} ของวัตถุ ตัวที่ i^{th} ด้วยสมการข้างล่าง เพื่อแปลงช่วงค่าที่เป็นไปได้ของตัวแปรให้อยู่ภาย ในช่วง $[0, 1]$

$$\frac{v' - \min'}{\max' - \min'} = \frac{v - \min}{\max - \min}$$

2110773-4 2/66

$$\frac{z_{if} - 0}{1 - 0} = \frac{r_{if} - 1}{M_f - 1}$$



17

Table 7.3 A sample data table containing variables of mixed type.

object identifier	test-1 (categorical)	test-2 (ordinal)	test-3 (ratio-scaled)
1	code-A	excellent	445
2	code-B	fair	22
3	code-C	good	164
4	code-A	excellent	1,210

log (test-3)

$$\begin{aligned} &\text{take log} \\ &2.65 \\ &1.34 \\ &2.21 \\ &3.08 \end{aligned}$$

fair 0
good 0.5
excellent 1

Dissimilarity Matrix of test-1

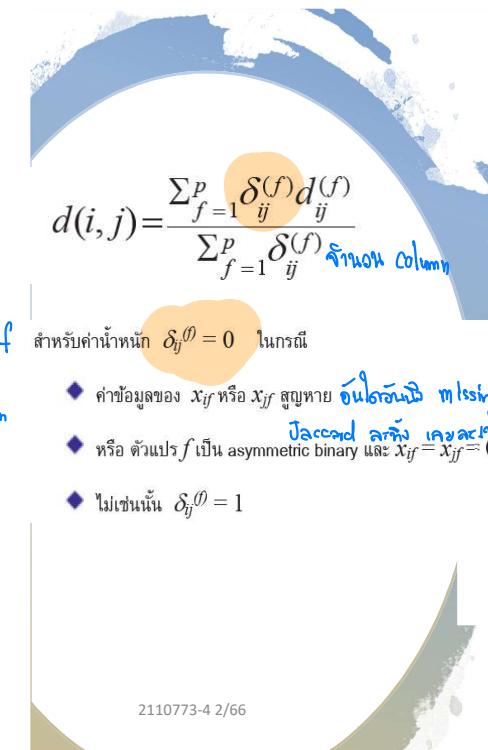
$$\begin{aligned} z_{if} &= \frac{r_{if} - 1}{M_f - 1} = \frac{1 - 0}{3 - 1} = 0 \\ d(2, 1) &= 0 \\ d(3, 1) &= 0 \\ d(3, 2) &= 0 \\ d(4, 1) &= 0 \\ d(4, 2) &= 0 \\ d(4, 3) &= 0 \end{aligned}$$

$$\frac{2-1}{3-1} = 0.5$$

$$\frac{3-1}{3-1} = 1$$

2110773-4 2/66

19



2110773-4 2/66

Variables of Mixed Types Attribute ที่มีหลาย data type

- ถ้า f เป็นตัวแปร nominal หรือ binary จะได้รับ $d_{ij}^{(f)} = 0$ หาก $x_{if} = x_{jf}$, หรือในกรณี $d_{ij}^{(f)} = 1$ สำหรับค่าใดๆ ก็ได้ $\delta_{ij}^{(f)} = 0$ ในกรณี
- ถ้า f เป็นตัวแปร interval-scaled ให้ใช้ระยะห่างของ mol ใช้ $d_{ij}^{(f)} = \frac{|x_{if} - x_{jf}|}{\max_h x_{hf} - \min_h x_{hf}}$ หรืออยู่ในช่วง $0-1$
- ถ้า f เป็นตัวแปร ordinal หรือ ratio-scaled
 - คำนวณค่าอันดับ r_{if} และ normalize ด้วยค่า $Z_{if} = \frac{r_{if} - 1}{M_f - 1}$
 - มองว่าค่า Z_{if} เป็นตัวแปร interval-scaled

18

we can use the dissimilarity matrices obtained for test-1 and test-2 later when we compute Equation (7.15). First, however, we need to complete some work for test-3 (which is ratio-scaled). We have already applied a logarithmic transformation to its values. Based on the transformed values of 2.65, 1.34, 2.21, and 3.08 obtained for the objects 1 to 4, respectively, we let $\max_h = 3.08$ and $\min_h = 1.34$. We then normalize the values in the dissimilarity matrix obtained in Example 7.5 by dividing each one by $(3.08 - 1.34) = 1.74$. This results in the following dissimilarity matrix for test-3:

$$\begin{bmatrix} 0 & 0.75 & 0 & 0 \\ 0.75 & 0 & 0.50 & 0 \\ 0 & 0.50 & 0 & 0 \\ 0.25 & 0.25 & 0 & 0 \end{bmatrix}$$

$\frac{1.34}{1.74} = 0.75$ $\frac{3.08}{1.74} = 1.74$ $\frac{1.34}{1.74} = 0.75$ $\frac{3.08}{1.74} = 1.74$ - normalize

We can now use the dissimilarity matrices for the three variables in our computation of Equation (7.15). For example, we get $d(2, 1) = \frac{1(1) + 1(1) + 1(0.75)}{3} = 0.92$. The resulting dissimilarity matrix obtained for the data described by the three variables of mixed types is:

$$\begin{bmatrix} 0 & 0.92 & 0.58 & 0.08 \\ 0.92 & 0 & 0.67 & 1.00 \\ 0.58 & 0.67 & 0 & 0.67 \\ 0.08 & 1.00 & 0.67 & 0 \end{bmatrix}$$

$$d(i, j) = \frac{\sum_{f=1}^p \delta_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^p \delta_{ij}^{(f)}},$$

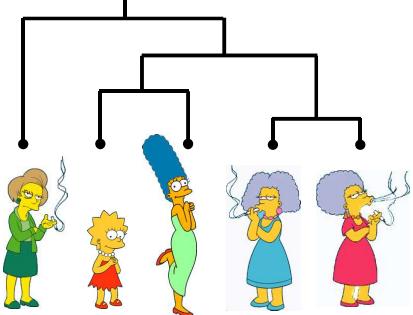
2110773-4 2/66

20

Two types of clustering

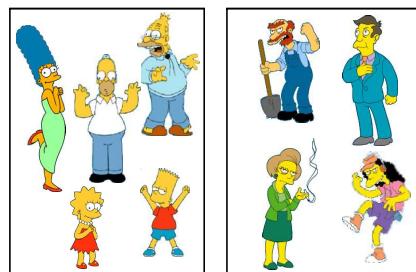
- **Partitional algorithms:** Construct various partitions and then evaluate them by some criterion
- **Hierarchical algorithms:** Create a hierarchical decomposition of the set of objects using some criterion

Hierarchical



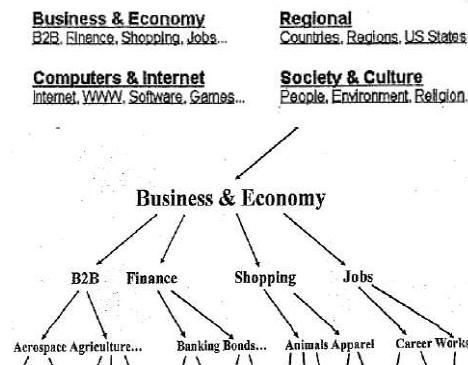
2110773-4 2/66

Partitional



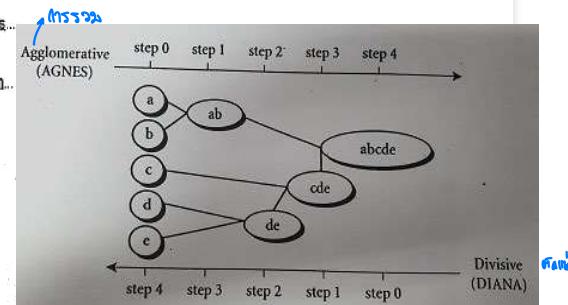
21

Top down vs. Bottom up

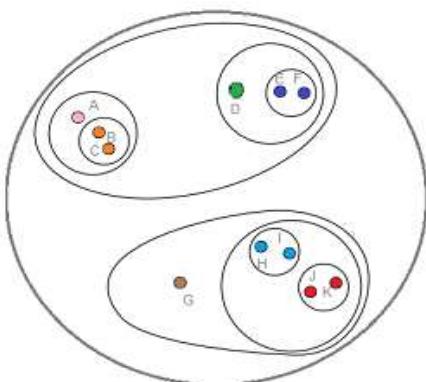


ตัวอย่างคำดำเนินการจัดกลุ่มเว็บไซต์เดรกทอรี

2110773-4 2/66



22



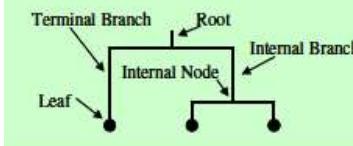
A dendrogram (right) representing nested clusters (left).

Dendrogram

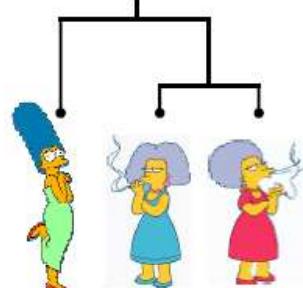
- [tree diagram](#) showing hierarchical clustering — relationships between similar sets of data.
- A tree that shows clustering process.

2110773-4 2/66

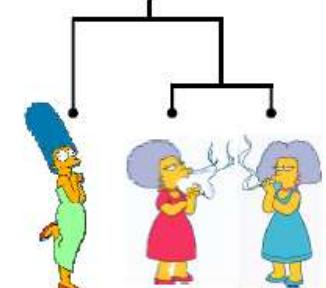
Parts of a Dendrogram



The similarity between two objects in a dendrogram is represented as the height of the lowest internal node they share.



21

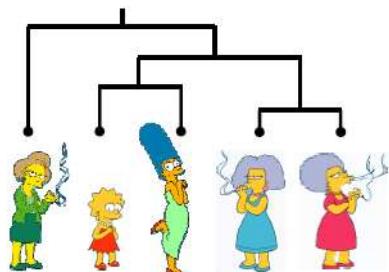


24

(How-to) Hierarchical Clustering

The number of dendograms with n leafs = $(2n - 3)! / [(2^{(n-2)}) (n - 2)!]$

Number of Leafs	Number of Possible Dendograms
2	1
3	3
4	15
5	105
...	...
10	34,459,425



2110773-4 2/66

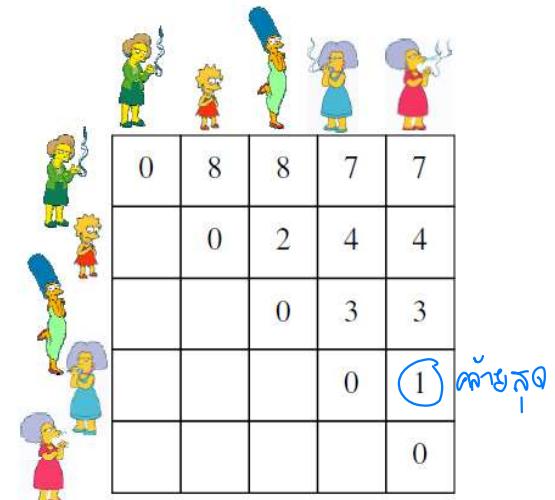
25

Since we cannot test all possible trees we will have to heuristic search of all possible trees. We could do this..

Bottom-Up (agglomerative): Starting with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.

Top-Down (divisive): Starting with all the data in a single cluster, consider every possible way to divide the cluster into two. Choose the best division and recursively operate on both sides.

We begin with a distance matrix which contains the distances between every pair of objects in our database.



26

$$D(\text{Marge}, \text{Homer}) = 8$$

$$D(\text{Bart}, \text{Lisa}) = 1$$

2110773-4 2/66

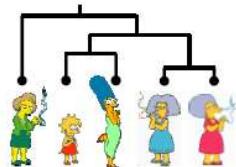
නැංවාක්ෂණීය

Bottom-Up (agglomerative):
Starting with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.

Consider all possible merges...

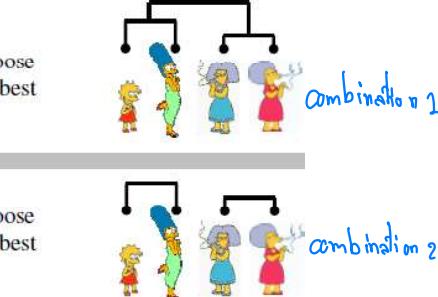


Choose the best



Consider all possible merges...

Choose the best



Consider all possible merges...

Choose the best



2110773-4 2/66

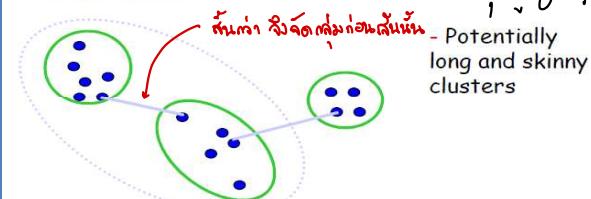
27

- **Single-linkage:** similarity of the closest pair. This can cause premature merging of groups with close pairs, even if those groups are quite dissimilar overall.

2110773-4 2/66

Similarity criteria: Single Link

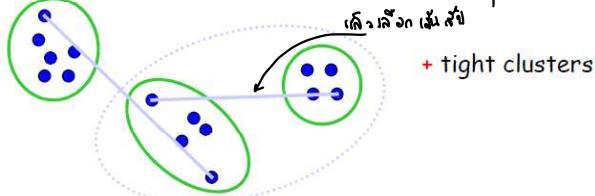
- cluster similarity = similarity of two **most similar members** member කිහිපෙන්ගේ group ගණ



28

Hierarchical: Complete Link

- cluster similarity = similarity of two least similar members



• **Complete linkage:**
similarity of the furthest pair. One drawback is that **outliers** can cause merging of close groups later than is optimal.

పాలోయిడ్ ఒట్లీ
గెంచికి సెంప్రెషన్ వైఫ్ లోవ్
ఎగ్జిస్టెంస్ క్లుబ్ గ్రాం

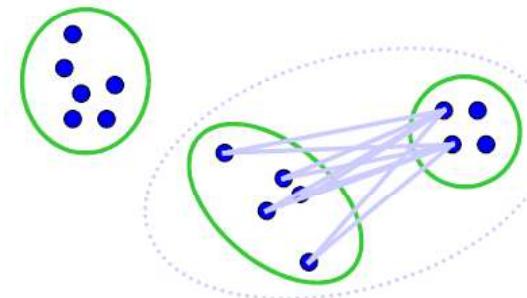
2110773-4 2/66

29

30

Hierarchical: Average Link

- cluster similarity = **average** similarity of all pairs



the most widely used similarity measure

Robust against noise

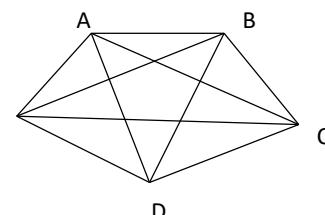
- similarity between groups.

An example to show different Links

Single link

- Merge the nearest clusters measured by the shortest edge between the two
- $((A\ B)\ (C\ D))\ E$

	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0



Complete link

- Merge the nearest clusters measured by the longest edge between the two
- $((A\ B)\ E)\ (C\ D)$

Average link

- Merge the nearest clusters measured by the average edge length between the two
- $((A\ B)\ (C\ D))\ E$

2110773-4 2/66

31

32

Example Merge w/ Complete Linkage (1)

samples	A	B	C	D	E	F	G
A	0	0.5000	0.4286	1.0000	0.2500	0.6250	0.3750
B	0.5000	0	0.7143	0.8333	0.6667	0.2000	0.7778
C	0.4286	0.7143	0	1.0000	0.4286	0.6667	0.3333
D	1.0000	0.8333	0.4286	0	1.0000	0.8000	0.8571
E	0.2500	0.6667	0.4286	0.2500	0	0.7778	0.3750
F	0.6250	0.2000	0.6667	0.8000	0.7778	0	0.7500
G	0.3750	0.7778	0.3333	0.8571	0.3750	0.7500	0

① ఇంధనంతరం

② మార్గాన్ని లేదా మార్గాన్ని

$A \rightarrow B$ ఏదు $A \rightarrow F$

0.5
0.6250

samples	A	(B,F)	C	D	E	G
A	0	0.6250	0.4286	1.0000	0.2500	0.3750
(B,F)	0.6250	0	0.7143	0.8333	0.7778	0.7778
C	0.4286	0.7143	0	1.0000	0.4286	0.3333
D	1.0000	0.8333	0.4286	0	1.0000	0.8571
E	0.2500	0.7778	0.4286	0.2500	0	0.3750
G	0.3750	0.7778	0.3333	0.8571	0.3750	0

పంచ స్థానించి ఉండవచ్చి

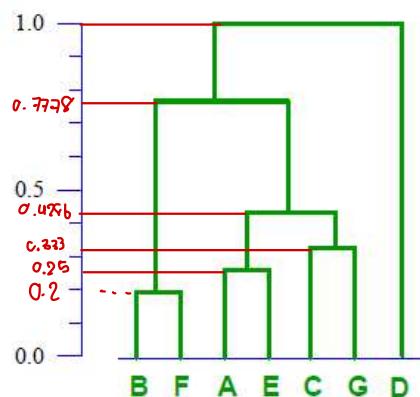
Example Merge w/ Complete Linkage (2)

samples	(A,E)	(B,F)	C	D	G
(A,E)	0	0.7778	0.4286	1.0000	0.3750
(B,F)	0.7778	0	0.7143	0.8333	0.7778
C	0.4286	0.7143	0	1.0000	0.3333
D	1.0000	0.8333	1.0000	0	0.8571
G	0.3750	0.7778	0.3333	0.8571	0

samples	(A,E)	(B,F)	(C,G)	D
(A,E)	0	0.7778	0.4286	1.0000
(B,F)	0.7778	0	0.7778	0.8333
(C,G)	0.4286	0.7778	0	1.0000
D	1.0000	0.8333	1.0000	0

samples	(A,E,C,G)	(B,F)	D
(A,E,C,G)	0	0.7778	1.0000
(B,F)	0.7778	0	0.8333
D	1.0000	0.8333	0

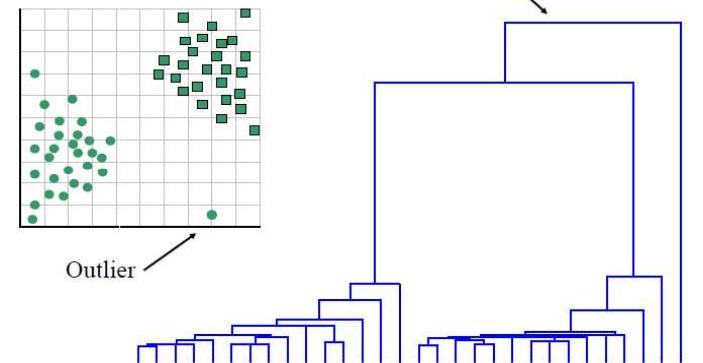
2110773-4 2/66



33

Dendrogram for outlier detection

The single isolated branch is suggestive of a data point that is very different to all others



34

Summary of hierachal clustering

- No need to specify the number of clusters in advance. *ກຳທັງຈະບຸວ່າຕ້ອງກຳລົງ*
- Hierarchical structure maps nicely onto human intuition for some domains
- They do not scale well: time complexity of at least $O(n^2)$, where n is the number of total objects.
- Like any heuristic search algorithms, local optima are a problem.
- Interpretation of results is (very) subjective.
- Only a tree is returned

2110773-4 2/66

Partitional Clustering *K - Mean*

- Nonhierarchical, each instance is placed in exactly one of K nonoverlapping clusters.
- Since only one set of clusters is output, the user normally has to input the desired number of clusters K.



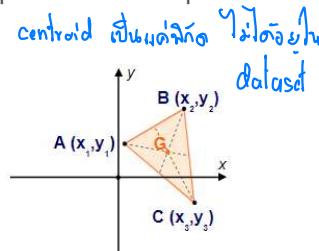
35

2110773-4 2/66

36

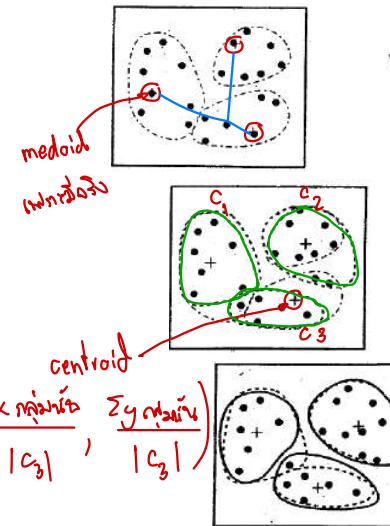
Centroid vs. Medoid

- Centroid, also called a geometric center, is the **center of mass of an object** (geometric balanced point)
- In general, a centroid is the *arithmetic mean* of all the points in the shape.



- Medoids are representative objects of a data set or a cluster within a data set whose sum of dissimilarities to all the objects in the cluster is minimal.
- Medoids are similar in concept to centroids, but medoids are always restricted to be members of the data set.

medoid ต้องมาจาก dataset



Algorithm: The k -Means algorithm for partitioning based on the mean value of objects in the cluster.

Input: The number of clusters k and a database containing n objects.

Output: A set of k clusters that minimizes the squared-error criterion.

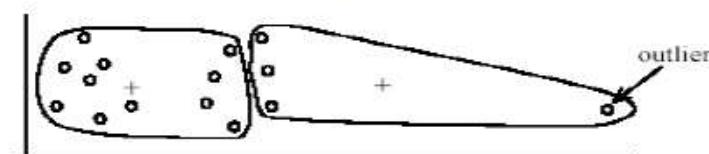
Method: คัดเลือก k ตัวอย่าง

- (1) arbitrarily choose k objects as the initial cluster centers;
- (2) repeat
- (3) (re)assign each object to the cluster to which the object is the most similar, based on the mean value of the objects in the cluster;
- (4) update the cluster means i.e., calculate the mean value of the objects for each cluster;
- (5) until no change;

K-means Summary

- Strength**
 - Simple, easy to implement and debug
 - Intuitive objective function: optimizes intra-cluster similarity
 - Relatively efficient: $O(tkn)$, where n is # objects, k is # clusters, and t is # iterations. Normally, $k, t \ll n$.
- Weakness**
 - Applicable only when *mean* is defined, then what about categorical data? กรณี mean ใช้ categorical ยังไง
 - Often terminates at a *local optimum*. Initialization is important.
 - Need to specify K , the *number* of clusters, in advance
 - Unable to handle noisy data and *outliers* จุด outlier ไม่เข้า
 - Not suitable to discover clusters with *non-convex shapes* รูปหินฟัน, รูปหัวใจ

Weaknesses of k-means: Problems with outliers



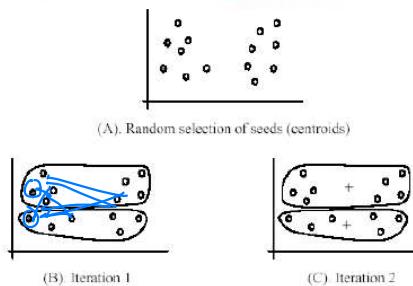
(A): Undesirable clusters



(B): Ideal clusters

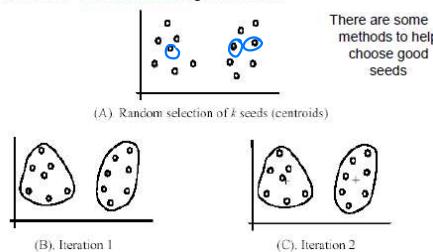
Weaknesses of k-means (cont ...)

- The algorithm is sensitive to **initial seeds**.



นี่ random ตั้งค่าเริ่มต้นไปต่อไม่ถูก

- If we use **different seeds**: good results



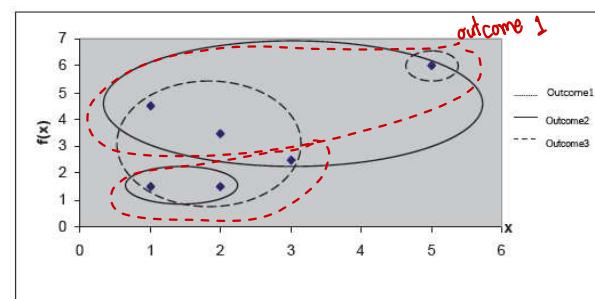
2110773-4 2/66

41

คุณภาพ
การจัดกลุ่ม

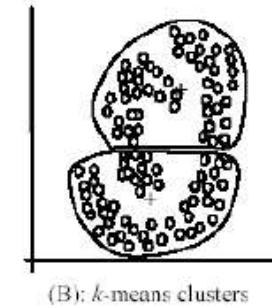
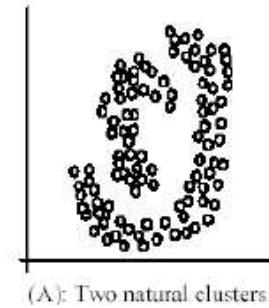
43

วัตถุ	x	y
1	1.0	1.5
2	1.0	4.5
3	2.0	1.5
4	2.0	3.5
5	3.0	2.5
6	5.0	6.0



Weaknesses of k-means (cont ...)

- The k -means algorithm is not suitable for discovering clusters that are not hyper-ellipsoids (or hyper-spheres).



2110773-4 2/66

42

2110773-4 2/66

Outcome	Cluster Centers	Cluster Points	Squared Error
1	(2.67, 4.67) (2.00, 1.83)	2, 4, 6 1, 3, 5	14.50
2	(1.5, 1.5) (2.75, 4.125)	1, 3 2, 4, 5, 6	15.94
3	(1.8, 2.7) (5.6)	1, 2, 3, 4, 5 6	9.60 high intra class similarity

$$E = \sum_{i=1}^k \sum_{p \in C_i} |p - \bar{m}_i|^2$$

ตัวนี้ดูเหมือนกัน

โดย ค่า E เป็นผลรวมค่าความผิดพลาดยกกำลังสองของวัตถุทั้งหมด

p คือ ค่าพิกัดในปริภูมิหลายมิติของวัตถุ

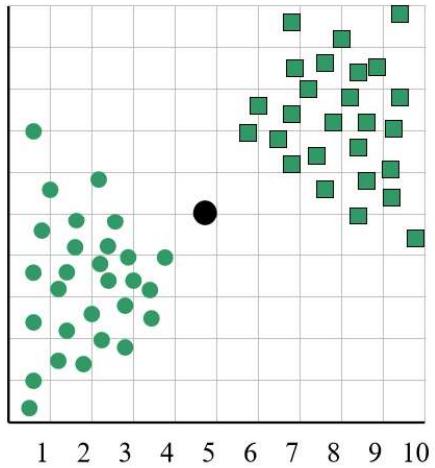
\bar{m}_i คือ ค่าพิกัดเฉลี่ยในปริภูมิหลายมิติของกลุ่ม C_i

2110773-4 2/66

44

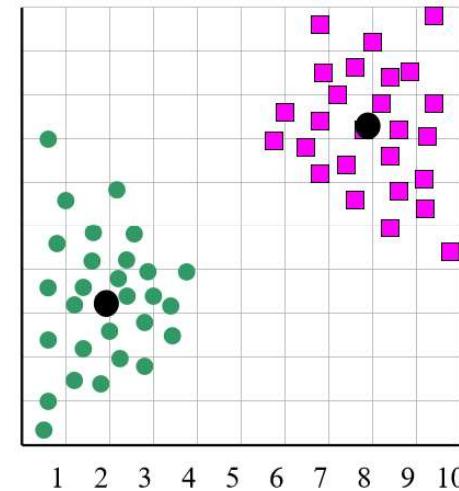
How can we tell the *right* number of clusters?

K=1



How can we tell the *right* number of clusters?

K=2



2110773-4 2/66

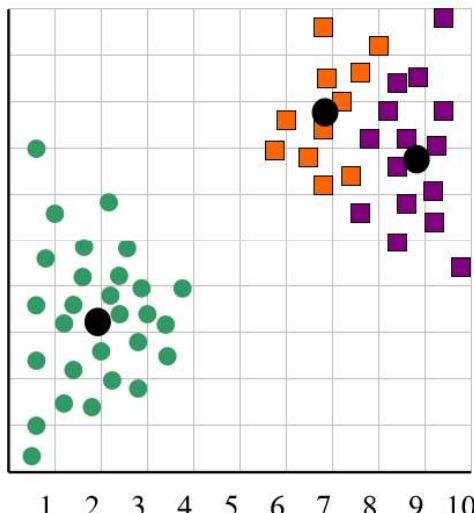
45

2110773-4 2/66

46

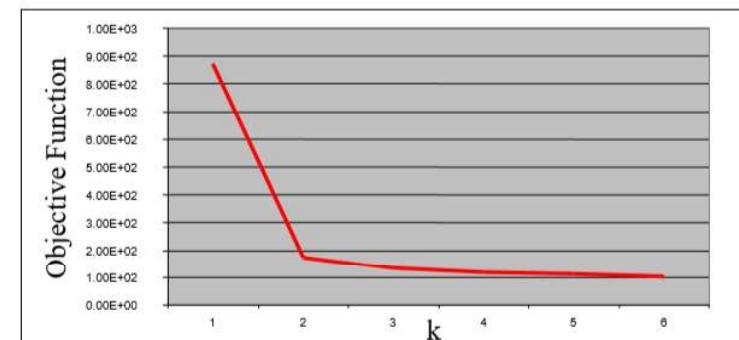
How can we tell the *right* number of clusters?

K=3



ก่อซักกวนอุนฟง: SSE
We can plot the objective function values for k equals 1 to 6...

The abrupt change at k = 2, is highly suggestive of two clusters in the data. This technique for determining the number of clusters is known as “knee finding” or “elbow finding”.



Note that the results are not always as clear cut as in this toy example

2110773-4 2/66

47

2110773-4 2/66

48