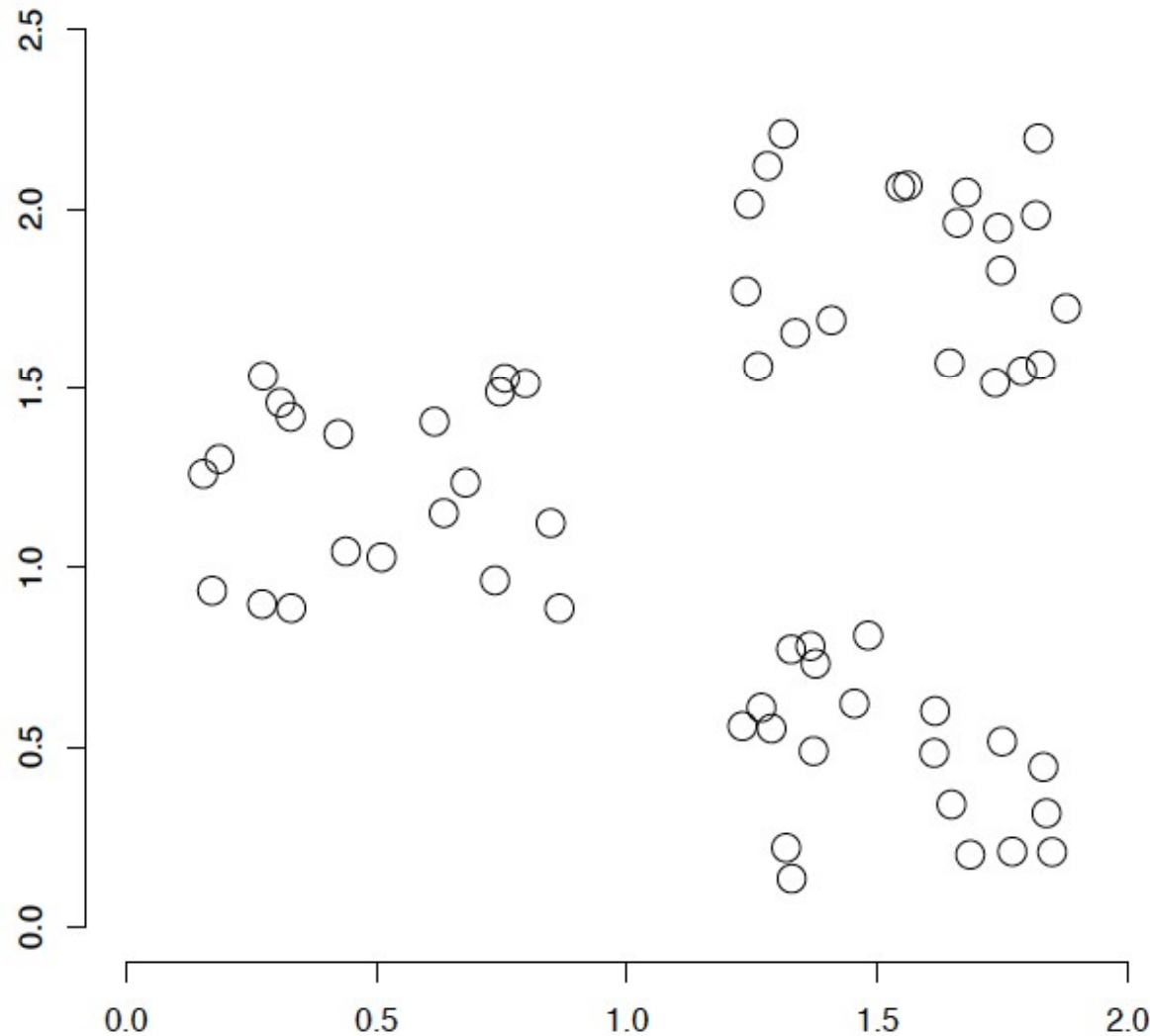


What is Cluster Analysis?

- ▶ **Cluster:** a collection of data objects
 - ▶ Similar to one another within the same cluster
 - ▶ Dissimilar to the objects in other clusters
- ▶ **Cluster analysis**
 - ▶ Finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters
- ▶ **Unsupervised learning:** no predefined classes
- ▶ Typical applications
 - ▶ As a **stand-alone tool** to get insight into data distribution
 - ▶ As a **preprocessing step** for other algorithms

A data set with clear cluster structure



► How would you design an algorithm for finding the three clusters in this case?

Examples of Clustering Applications

- ▶ Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs
- ▶ Land use: Identification of areas of similar land use in an earth observation database
- ▶ Insurance: Identifying groups of motor insurance policy holders with a high average claim cost
- ▶ City-planning: Identifying groups of houses according to their house type, value, and geographical location

Measure the Quality of Clustering

- ▶ Dissimilarity/Similarity metric: Similarity is expressed in terms of a distance function, typically metric: $d(i, j)$
- ▶ There is a separate “quality” function that measures the “goodness” of a cluster.
- ▶ The definitions of distance functions are usually very different for interval-scaled, boolean, categorical, ordinal ratio, and vector variables.
- ▶ Weights should be associated with different variables based on applications and data semantics.
- ▶ It is hard to define “similar enough” or “good enough”
 - ▶ the answer is typically highly subjective.

Requirements of Clustering in Data Mining

- ▶ Scalability
- ▶ Ability to deal with different types of attributes
- ▶ Ability to handle dynamic data
- ▶ Discovery of clusters with arbitrary shape
- ▶ Minimal requirements for domain knowledge to determine input parameters
- ▶ Able to deal with noise and outliers
- ▶ High dimensionality

Data Structures

► Data matrix

- Two mode matrix : it has two kinds of entities i.e. dissimilarity
- n objects and p attributes

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

► Dissimilarity matrix

- One mode matrix : it has only one kind of entity i.e. dissimilarity
- $d(i,j)$ = dissimilarity measure

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

► Similarity matrix

- $\text{sim}(i,j) = 1 - d(i,j)$

Type of data in clustering analysis

- ▶ Interval-scaled variables
- ▶ Nominal, ordinal
- ▶ Variables of mixed types



Interval-valued variables

- ▶ Standardize data

- ▶ Calculate the mean absolute deviation:

$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

where

$$m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf}).$$

- ▶ Calculate the standardized measurement (z-score)

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

- ▶ Using mean absolute deviation is more robust than using standard deviation

Similarity and Dissimilarity Between Objects

(Discussed at time of preprocessing)

- ▶ Distances are normally used to measure the similarity or dissimilarity between two data objects

- ▶ 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

- ▶ If $q = 1$, d is Manhattan distance

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

Similarity and Dissimilarity Between Objects

(Discussed at time of preprocessing)

- Distances are normally used to measure the similarity or dissimilarity between two data objects

- 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

- If $q = 1$, d is Manhattan distance

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

Similarity and Dissimilarity Between Objects

(Discussed at time of preprocessing)

- ▶ If $q = 2$, d is Euclidean distance:

$$d(i,j) = \sqrt{(|x_{i_1} - x_{j_1}|^2 + |x_{i_2} - x_{j_2}|^2 + \dots + |x_{i_p} - x_{j_p}|^2)}$$

- ▶ Properties

- ▶ $d(i,j) \geq 0$
 - ▶ $d(i,i) = 0$
 - ▶ $d(i,j) = d(j,i)$
 - ▶ $d(i,j) \leq d(i,k) + d(k,j)$
- ▶ Also, one can use weighted distance, parametric Pearson product moment correlation, or other dissimilarity measures

Dissimilarity for Nominal Variables

- Can take 2 or more states, e.g., red, yellow, blue, green (generalization of a binary attribute)
- Method 1: Simple matching

- m : # of matches, p : total # of variables

$$d(i, j) = \frac{p - m}{p}$$

Object Identifier	test-1 (nominal)
1	code A
2	code B
3	code C
4	code A

$$\begin{bmatrix} 0 & & & \\ d(2, 1) & 0 & & \\ d(3, 1) & d(3, 2) & 0 & \\ d(4, 1) & d(4, 2) & d(4, 3) & 0 \end{bmatrix}.$$

Since here we have one nominal attribute, *test-1*, we set $p = 1$ in Eq. (2.11) so that $d(i, j)$ evaluates to 0 if objects i and j match, and 1 if the objects differ. Thus, we get

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}.$$

From this, we see that all objects are dissimilar except objects 1 and 4 (i.e., $d(4, 1) = 0$).

Dissimilarity for Binary Variables

- A contingency table for binary data

		Object j		
		1	0	sum
Object i	1	q	r	$q+r$
	0	s	t	$s+t$
sum		$q+s$	$r+t$	p

- Distance measure for symmetric binary variables:

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

- Distance measure for asymmetric binary variables:

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

- Jaccard coefficient (*similarity* measure for *asymmetric* binary variables):

$$sim_{Jaccard}(i, j) = \frac{q}{q + r + s}$$

Note: $p(q+r+s+t)$ is total no of attributes.

Dissimilarity between Binary Variables

► Example

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

- gender is a symmetric attribute
- the remaining attributes are asymmetric binary
- let the values Y and P be set to 1, and the value N be set to 0

$$d (jack , mary) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d (jack , jim) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d (jim , mary) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$

Dissimilarity between Ordinal Variables

- An ordinal variable can be discrete or continuous
- Order is important, e.g., rank
- Can be treated like interval-scaled
 - replace x_{if} by their rank $r_{if} \in \{1, \dots, M_f\}$
 - map the range of each variable onto $[0, 1]$ by replacing i -th object in the f -th variable by

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

- compute the dissimilarity using methods for interval-scaled variables

A Sample Data Table Containing Attributes of Mixed Type

Object Identifier	test-1 (nominal)	test-2 (ordinal)	test-3 (numeric)
1	code A	excellent	45
2	code B	fair	22
3	code C	good	64
4	code A	excellent	28

$$\begin{bmatrix} 0 & & & \\ 1.0 & 0 & & \\ 0.5 & 0.5 & 0 & \\ 0 & 1.0 & 0.5 & 0 \end{bmatrix}^*$$

Ordinal attribute: test-2,

Three states for test-2: fair, good, and excellent, that is, $M_f = 3$.

Step 1, if we replace each value for test-2 by its rank, the four objects are assigned the ranks 3, 1, 2, and 3, respectively.

Step 2 normalizes the ranking by mapping rank 1 to 0.0, rank 2 to 0.5, and rank 3 to 1.0.

Step 3, we can use, say, the Euclidean distance which results in the given dissimilarity matrix:

Dissimilarity for Mixed Variable

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

where the indicator $\delta_{ij}^{(f)} = 0$ if either (1) x_{if} or x_{jf} is missing (i.e., there is no measurement of attribute f for object i or object j), or (2) $x_{if} = x_{jf} = 0$ and attribute f is asymmetric binary; otherwise, $\delta_{ij}^{(f)} = 1$. The contribution of attribute f to the dissimilarity between i and j (i.e., $d_{ij}^{(f)}$) is computed dependent on its type:

- If f is numeric: $d_{ij}^{(f)} = \frac{|x_{if} - x_{jf}|}{\max_h x_{hf} - \min_h x_{hf}}$, where h runs over all nonmissing objects for attribute f .
- If f is nominal or binary: $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$; otherwise, $d_{ij}^{(f)} = 1$.
- If f is ordinal: compute the ranks r_{if} and $z_{if} = \frac{r_{if} - 1}{M_f - 1}$, and treat z_{if} as numeric.

Example: Dissimilarity for Mixed Variable

A Sample Data Table Containing Attributes of Mixed Type

Object Identifier	test-1 (nominal)	test-2 (ordinal)	test-3 (numeric)
1	code A	excellent	45
2	code B	fair	22
3	code C	good	64
4	code A	excellent	28

Test 1 =>
$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

Test 2 =>
$$\begin{bmatrix} 0 & & & \\ 1.0 & 0 & & \\ 0.5 & 0.5 & 0 & \\ 0 & 1.0 & 0.5 & 0 \end{bmatrix}$$

Test 3 =>
$$\begin{bmatrix} 0 & & & \\ 0.55 & 0 & & \\ 0.45 & 1.00 & 0 & \\ 0.40 & 0.14 & 0.86 & 0 \end{bmatrix}$$

$D(3,1) \Rightarrow d(3,1) = \frac{1(1)+1(0.50)+1(0.45)}{3} = 0.65$, indicator $\delta_{ij}^{(q)} = 1$ for each of the three attributes

Overall Dissimilarity matrix =>
$$\begin{bmatrix} 0 & & & \\ 0.85 & 0 & & \\ 0.65 & 0.83 & 0 & \\ 0.13 & 0.71 & 0.79 & 0 \end{bmatrix}$$

Cosine Similarity between documents

- A **document** can be represented by thousands of attributes, each recording the *frequency* of a particular word (such as keywords) or phrase in the document.

Document	team	coach	hockey	baseball	soccer	penalty	score	win	loss	season
Document1	5	0	3	0	2	0	0	2	0	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	0
Document4	0	1	0	0	1	2	2	0	3	0

- Other vector objects: gene features in micro-arrays, ...
- Applications: information retrieval, biologic taxonomy, gene feature mapping, ...
- Cosine measure: If d_1 and d_2 are two vectors (e.g., term-frequency vectors), then

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

where \bullet indicates vector dot product, $||d||$: the length of vector d

Cosine Similarity between documents

- $\cos(d_1, d_2) = (d_1 \bullet d_2) / ||d_1|| ||d_2||$,
where \bullet indicates vector dot product, $||d||$: the length of vector d
- Ex: Find the similarity between documents 1 and 2.

$$d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$d_2 = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$$

$$d_1 \bullet d_2 = 5*3 + 0*0 + 3*2 + 0*0 + 2*1 + 0*1 + 0*1 + 2*1 + 0*0 + 0*1 = 25$$

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

$$||d_2|| = (3*3 + 0*0 + 2*2 + 0*0 + 1*1 + 1*1 + 0*0 + 1*1 + 0*0 + 1*1)^{0.5} = (17)^{0.5} = 4.12$$

$$\cos(d_1, d_2) = 0.94$$

Major Clustering Approaches

- ▶ Partitioning approach:

- ▶ Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors
- ▶ Typical methods: **k-means**, **k-medoids**, CLARANS

- ▶ Hierarchical approach:

- ▶ Create a hierarchical decomposition of the set of data (or objects) using some criterion
- ▶ Typical methods: **Diana**, **Agnes**, BIRCH, ROCK, CAMELEON

- ▶ Density-based approach:

- ▶ Based on connectivity and density functions
- ▶ Typical methods: **DBSCAN**, OPTICS, DenClue

Partitioning Algorithms: Basic Concept

- ▶ Partitioning method: Construct a partition of a database D of n objects into a set of k clusters, s.t., min sum of squared distance

$$\sum_{m=1}^k \sum_{t_{mi} \in K_m} (C_m - t_{mi})^2$$

- ▶ Given a k , find a partition of k clusters that optimizes the chosen partitioning criterion
 - ▶ Global optimal: exhaustively enumerate all partitions
 - ▶ Heuristic methods: k -means and k -medoids algorithms
 - ▶ k -means (MacQueen'67): Each cluster is represented by the center of the cluster
 - ▶ k -medoids or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster

Cluster Analysis

1. What is Cluster Analysis?
2. A Categorization of Major Clustering Methods
3. **Partitioning Methods**
4. Hierarchical Methods
5. Density-Based Methods
6. Summary

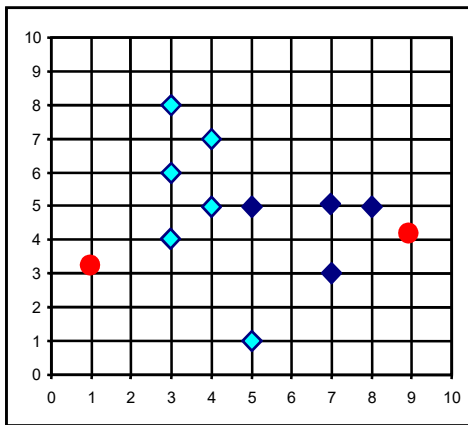


The *K-Means* Clustering Method

- ▶ Given k , the *k-means* algorithm is implemented in four steps:
 - ▶ Partition objects into k nonempty subsets
 - ▶ Compute seed points as the centroids of the clusters of the current partition (the centroid is the center, i.e., *mean point*, of the cluster)
 - ▶ Assign each object to the cluster with the nearest seed point
 - ▶ Go back to Step 2, stop when no more new assignment to cluster(centroids remains unchanged / some threshold is reached)

The *K*-Means Clustering Method

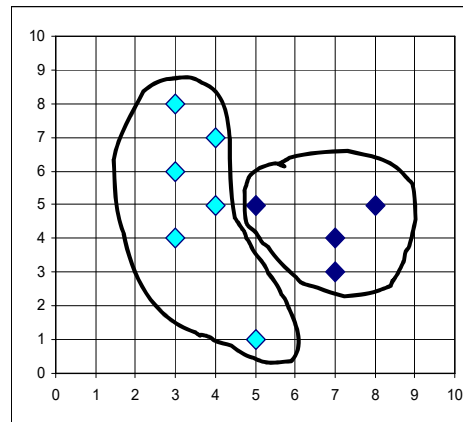
► Example



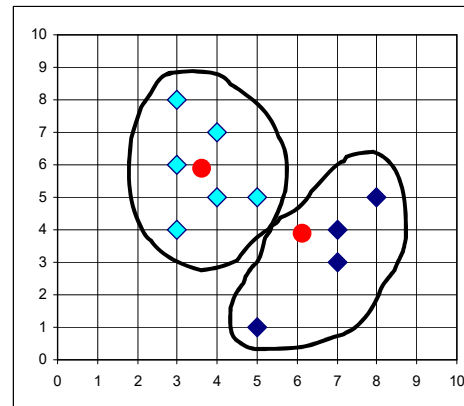
$K=2$

Arbitrarily choose K object as initial cluster center

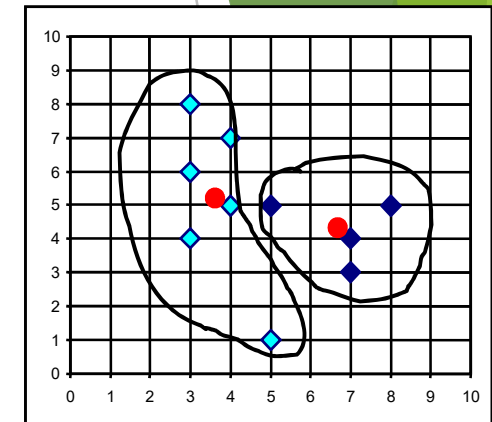
Assign each object to most similar center



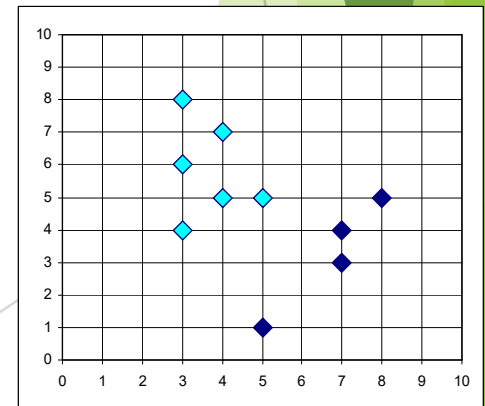
reassign



Update the cluster means



reassign



Update the cluster means

Example



Comments on the *K-Means* Method

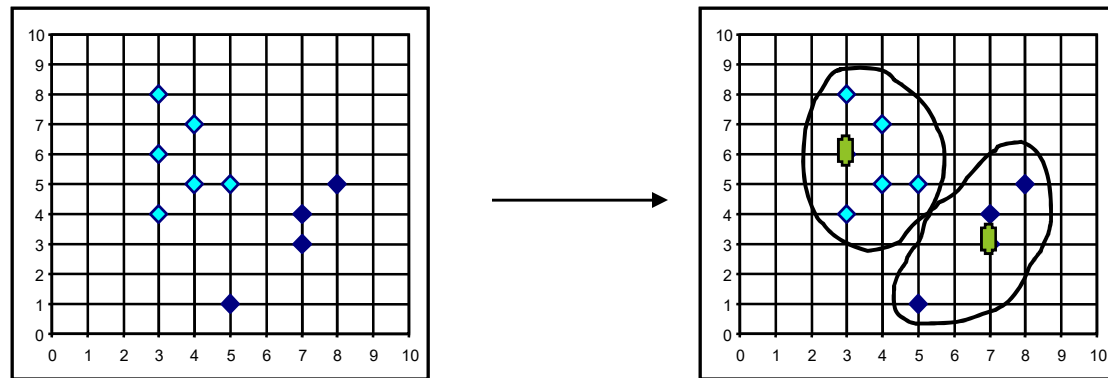
- ▶ Strength: *Relatively efficient*: $O(tkn)$, where n is # objects, k is # clusters, and t is # iterations. Normally, $k, t \ll n$.
 - ▶ Comparing: PAM: $O(k(n-k)^2)$
- ▶ Comment: Often terminates at a *local optimum*. The *global optimum* may be found using techniques such as *genetic algorithms*
- ▶ Weakness
 - ▶ Applicable only when *mean* is defined, then what about categorical data?
 - ▶ Need to specify k , the *number* of clusters, in advance
 - ▶ Unable to handle noisy data and *outliers*
 - ▶ Not suitable to discover clusters with *non-spherical shapes*

Variations of the *K-Means* Method

- ▶ A few variants of the *k-means* which differ in
 - ▶ Selection of the initial *k* means
 - ▶ Dissimilarity calculations
- ▶ Handling categorical data: *k-modes* (Huang'98)
 - ▶ Replacing means of clusters with modes
 - ▶ Using new dissimilarity measures to deal with categorical objects
 - ▶ Using a frequency-based method to update modes of clusters
 - ▶ A mixture of categorical and numerical data: *k-prototype* method

What Is the Problem of the K-Means Method?

- ▶ The **k-means algorithm is sensitive to outliers** !
 - ▶ Since an object with an extremely large value may substantially distort the distribution of the data.
- ▶ **K-Medoids**: Instead of taking the **mean** value of the object in a cluster as a reference point, **medoids** can be used, which is the **most centrally located object** in a cluster.



The *K-Medoids* Clustering Method

- ▶ Minimize the sensitivity of k-means to outliers
- ▶ Pick actual objects to represent clusters instead of mean values
- ▶ Each remaining object is clustered with the representative object (**Medoid**) to which is the most similar
- ▶ The algorithm minimizes the sum of the dissimilarities between each object and its corresponding reference point

$$E = \sum_{i=1}^k \sum_{p \in C_i} |p - o_i|$$

- **E**: the sum of absolute error for all objects in the data set
- **P**: the data point in the space representing an object
- **O_i**: is the representative object of cluster C_i

K-Medoids Algorithm(PAM)

PAM : Partitioning Around Medoids

► **Input**

- K: the number of clusters
- D: a data set containing n objects

► **Output:** A set of k clusters

► **Method:**

- (1) Arbitrary choose k objects from D as representative objects (seeds)
- (2) **Repeat**
- (3) Assign each remaining object to the cluster with the nearest representative object
- (4) For each representative object O_j
- (5) Randomly select a non representative object O_{random}
- (6) Compute the total cost S of swapping representative object O_j with O_{random}
- (7) if $S < 0$ then replace O_j with O_{random}
- (8) **Until** no change

Example

No of clusters are 2

Point	x-axis	y-axis
1	7	6
2	2	6
3	3	8
4	8	5
5	7	4
6	4	7
7	6	2
8	7	3
9	6	4
10	3	4

What Is the Problem with PAM?

- ▶ Pam is more robust than k-means in the presence of noise and outliers because a medoid is less influenced by outliers or other extreme values than a mean
- ▶ Pam works efficiently for small data sets but does not scale well for large data sets.

- ▶ $O(k(n-k)^2)$ for each iteration

where n is # of data, k is # of clusters

→ Sampling based method,

CLARA(Clustering LARge Applications)

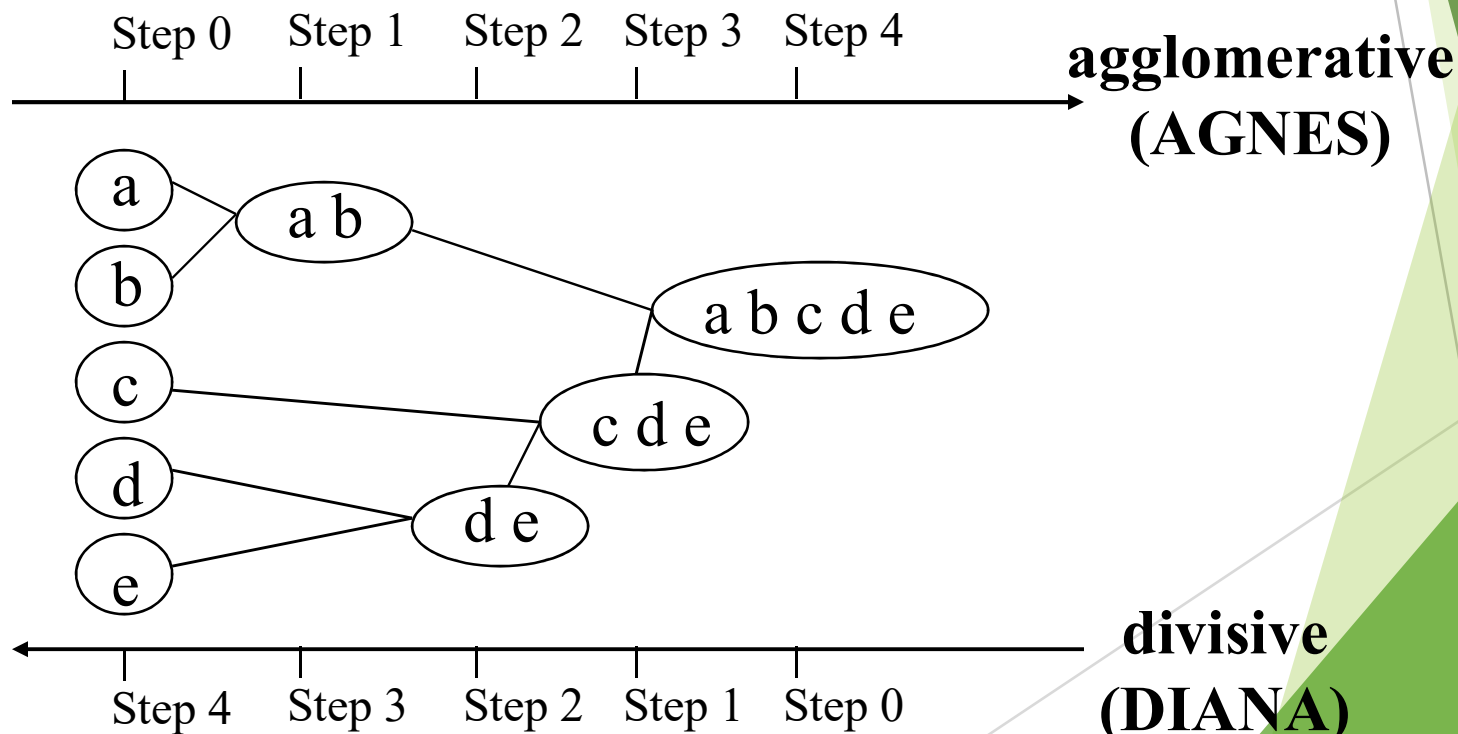
Cluster Analysis

1. What is Cluster Analysis?
2. A Categorization of Major Clustering Methods
3. Partitioning Methods
4. **Hierarchical Methods**
5. Density-Based Methods
6. Summary



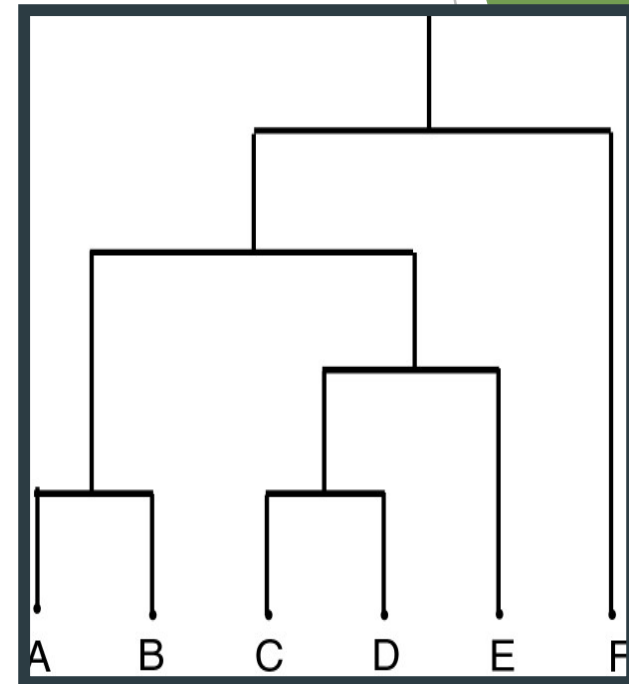
Hierarchical Clustering

- Use **distance matrix** as **clustering criteria**. This method does not require the number of clusters k as an input, but **needs a termination condition**
- Produces a dendrogram , hierarchical tree of clusters



Dendrogram

- ▶ **Dendrogram:** a tree data structure which illustrates hierarchical clustering techniques.
- ▶ Each level shows clusters for that level.
 - ▶ **Leaf** - individual clusters
 - ▶ **Root** - one cluster
- ▶ A cluster at level i is the union of its children clusters at level $i+1$.

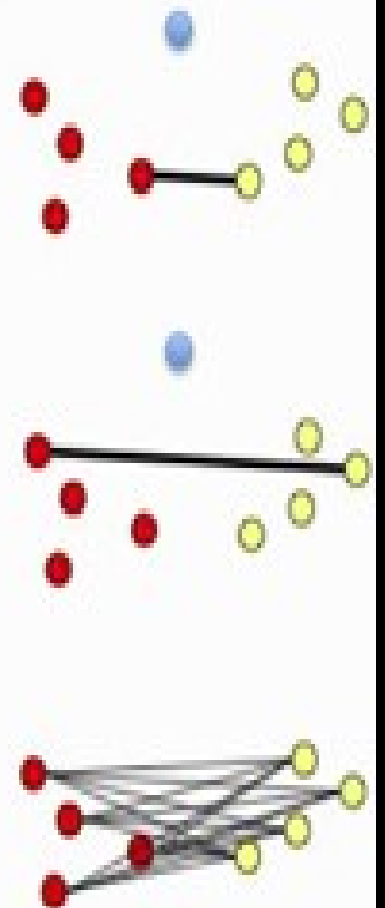


A gglomerative Hierarchical Clustering Algorithm

-
- 1: Compute the proximity matrix, if necessary.
 - 2: repeat
 - 3: Merge the closest two clusters.
 - 4: Update the proximity matrix to reflect the proximity between the new cluster and the original clusters.
 - 5: until Only one cluster remains.
-

Cluster distance measures

- Single link: $D(c_1, c_2) = \min_{x_1 \in c_1, x_2 \in c_2} D(x_1, x_2)$
 - distance between closest elements in clusters
 - produces long chains $a \rightarrow b \rightarrow c \rightarrow \dots \rightarrow z$
- Complete link: $D(c_1, c_2) = \max_{x_1 \in c_1, x_2 \in c_2} D(x_1, x_2)$
 - distance between farthest elements in clusters
 - forces "spherical" clusters with consistent "diameter"
- Average link: $D(c_1, c_2) = \frac{1}{|c_1|} \frac{1}{|c_2|} \sum_{x_1 \in c_1} \sum_{x_2 \in c_2} D(x_1, x_2)$
 - average of all pairwise distances
 - less affected by outliers



Agglomerative Clustering(AGNES)

Given a set of N items to be clustered, and an $N \times N$ distance (or similarity) matrix, the basic process of hierarchical clustering (defined by S.C. Johnson in 1967) is this:

- 1- Start by assigning each item to a cluster, so that if you have N items, you now have N clusters, each containing just one item. Let the distances (similarities) between the clusters the same as the distances (similarities) between the items they contain.
- 2- Find the closest (most similar) pair of clusters and merge them into a single cluster, so that now you have one cluster less.
- 3- Compute distances (similarities) between the new cluster and each of the old clusters.

Example (input is a set of points)

Step 1: calculate the distance matrix. Find the minimum value in distance matrix.

Dist	A	B	C	D	E	F
A	0.00	0.71	5.66	3.61	4.24	3.20
B	0.71	0.00	4.95	2.92	3.54	2.50
C	5.66	4.95	0.00	2.24	1.41	2.50
D	3.61	2.92	2.24	0.00	1.00	0.50
E	4.24	3.54	1.41	1.00	0.00	1.12
F	3.20	2.50	2.50	0.50	1.12	0.00

Step 2: Merge the points in one cluster and recalculate the distances.

Min Distance (Single Linkage)

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	?	4.24
B	0.71	0.00	4.95	?	3.54
C	5.66	4.95	0.00	?	1.41
D, F	?	?	?	0.00	?
E	4.24	3.54	1.41	?	0.00

- ▶ Using the input distance matrix, distance between cluster (D, F) and cluster A is computed as

$$d_{(D,F) \rightarrow A} = \min(d_{DA}, d_{FA}) = \min(3.61, 3.20) = 3.20$$

- ▶ Distance between cluster (D, F) and cluster B is

$$d_{(D,F) \rightarrow B} = \min(d_{DB}, d_{FB}) = \min(2.92, 2.50) = 2.50$$

- ▶ Similarly, distance between cluster (D, F) and cluster C is

$$d_{(D,F) \rightarrow C} = \min(d_{DC}, d_{FC}) = \min(2.24, 2.50) = 2.24$$

- ▶ Finally, distance between cluster E and cluster (D, F) is calculated as

- ▶
$$d_{E \rightarrow (D,F)} = \min(d_{ED}, d_{EF}) = \min(1.00, 1.12) = 1.00$$

- ▶ Then, the

- Step 3: update the matrix.

Min Distance (Single Linkage)

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	3.20	4.24
B	0.71	0.00	4.95	2.50	3.54
C	5.66	4.95	0.00	2.24	1.41
D, F	3.20	2.50	2.24	0.00	1.00
E	4.24	3.54	1.41	1.00	0.00

Step 4 : repeat step [1:3] until one cluster is made.

Dist	A,B	C	(D, F)	E
A,B	0	?	?	?
C	?	0	2.24	1.41
(D, F)	?	2.24	0	1.00
E	?	1.41	1.00	0

$$d_{C \rightarrow (A,B)} = \min(d_{CA}, d_{CB}) = \min(5.66, 4.95) = 4.95$$

$$d_{(D,F) \rightarrow (A,B)} = \min(d_{DA}, d_{DB}, d_{FA}, d_{FB}) = \min(3.61, 2.92, 3.20, 2.50) = 2.50$$

$$d_{E \rightarrow (A,B)} = \min(d_{EA}, d_{EB}) = \min(4.24, 3.54) = 3.54$$

Updated Matrix :

Min Distance (Single Linkage)

Dist	A,B	C	(D, F)	E
A,B	0	4.95	2.50	3.54
C	4.95	0	2.24	1.41
(D, F)	2.50	2.24	0	1.00
E	3.54	1.41	1.00	0

we can see that the closest distance between clusters happens between cluster E and (D, F) at distance 1.00. Thus, we cluster them together into cluster ((D, F), E).

Min Distance (Single Linkage)

Dist	(A,B)	C	(D, F), E
(A,B)	0.00	4.95	2.50
C	4.95	0.00	1.41
(D, F), E	2.50	1.41	0.00

$$d_{((D,F),E) \rightarrow (A,B)} = \min(d_{DA}, d_{DB}, d_{FA}, d_{FB}, d_{EA}, d_{EB}) = \min(3.61, 2.92, 3.20, 2.50, 4.24, 3.54) = 2.50$$

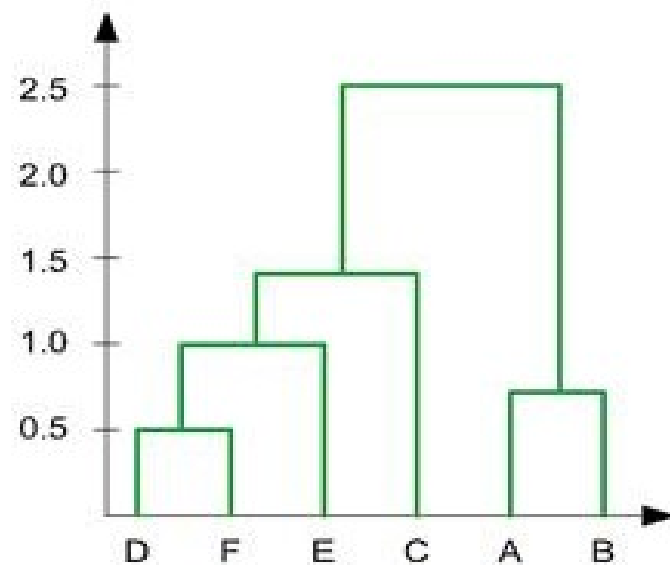
$$d_{((D,F),E) \rightarrow C} = \min(d_{DC}, d_{FC}, d_{EC}) = \min(2.24, 2.50, 1.41) = 1.41$$

$$d_{(((D,F),E),C) \rightarrow (A,B)} = \min(d_{DA}, d_{DB}, d_{FA}, d_{FB}, d_{EA}, d_{EB}, d_{CA}, d_{CB})$$

$$d_{(((D,F),E),C) \rightarrow (A,B)} = \min(3.61, 2.92, 3.20, 2.50, 4.24, 3.54, 5.66, 4.95) = 2.50$$

Min Distance (Single Linkage)

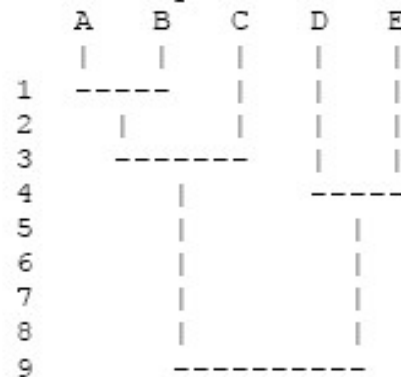
Dist	(A,B)	((D, F), E),C
(A,B)	0.00	2.50
((D, F), E),C	2.50	0.00



Complete Linkage : Example

	A	B	C	D	E
A	0	1	2	7	5
B		0	3	8	6
C			0	5	9
D				0	4
E					0

The result of complete-link clustering is the following dendrogram:



DIVISIVE Algorithms

- ▶ All the data is **part of one cluster initially** .
- ▶ Use a distance criterion to divide the cluster in two, and then subdivide the clusters until a stopping criterion is achieved.
- ▶ **Polythetic** - divide the data based on the values by all attributes (see the example)
- ▶ **Monothetic** - divide the data on the basis of the possession of a single specified attribute

Do this

Example (Polythetic)

- Use Average linkage distance if the the distance matrix is given.

	1	2	3	4	5	6	7
1	0						
2	10	0					
3	7	7	0				
4	30	23	21	0			
5	29	25	22	7	0		
6	38	34	31	10	11	0	
7	42	36	36	13	17	9	0

Polythetic Approach

	1	2	3	4	5	6	7	
1	0							$D(1, *) = 26.0$
2	10	0						$D(2, *) = 22.5$
3	7	7	0					$D(3, *) = 20.7$
4	30	23	21	0				$D(4, *) = 17.3$
5	29	25	22	7	0			$D(5, *) = 18.5$
6	38	34	31	10	11	0		$D(6, *) = 22.2$
7	42	36	36	13	17	9	0	$D(7, *) = 25.5$
$A = \{1\}$								
$B = \{2, 3, 4, 5, 6, 7\}$								

	1	2	3	4	5	6	7		
1	0							$D(2, A) = 10$	$D(2, B) = 25.0$
2	10	0						$D(3, A) = 7$	$D(3, B) = 23.4$
3	7	7	0					$D(4, A) = 30$	$D(4, B) = 14.8$
4	30	23	21	0				$D(5, A) = 29$	$D(5, B) = 16.4$
5	29	25	22	7	0			$D(6, A) = 38$	$D(6, B) = 19.0$
6	38	34	31	10	11	0		$D(7, A) = 42$	$D(7, B) = 22.2$
7	42	36	36	13	17	9	0		
$A = \{1\}$									
$B = \{2, 3, 4, 5, 6, 7\}$									

Initial Cluster
 $A = \{1\}$
 and
 $B = \{2, 3, 4, 5, 6, 7\}$

	1	2	3	4	5	6	7			
1	0							$D(2, A) = 10$	$D(2, B) = 25.0$	$\Delta_2 = 15.0$
2	10	0						$D(3, A) = 7$	$D(3, B) = 23.4$	$\Delta_3 = 16.4$
3	7	7	0					$D(4, A) = 30$	$D(4, B) = 14.8$	$\Delta_4 = -15.2$
4	30	23	21	0				$D(5, A) = 29$	$D(5, B) = 16.4$	$\Delta_5 = -12.6$
5	29	25	22	7	0			$D(6, A) = 38$	$D(6, B) = 19.0$	$\Delta_6 = -19.0$
6	38	34	31	10	11	0		$D(7, A) = 42$	$D(7, B) = 22.2$	$\Delta_7 = -19.8$
7	42	36	36	13	17	9	0			

$A = \{1, 3\}$	$D(7, A) = 42$	$D(7, B) = 22.2$	$\Delta_7 = -19.8$
$B = \{2, \text{✗} 4, 5, 6, 7\}$			

Updated Cluster
 $A = \{1, 3\}$
 and
 $B = \{2, 4, 5, 6, 7\}$

	1	2	3	4	5	6	7			
1	0							$D(2, A) = 10$	$D(2, B) = 25.0$	$\Delta_2 = 15.0$
2	10	0						$D(3, A) = 7$	$D(3, B) = 23.4$	$\Delta_3 = 16.4$
3	7	7	0					$D(4, A) = 30$	$D(4, B) = 14.8$	$\Delta_4 = -15.2$
4	30	23	21	0				$D(5, A) = 29$	$D(5, B) = 16.4$	$\Delta_5 = -12.6$
5	29	25	22	7	0			$D(6, A) = 38$	$D(6, B) = 19.0$	$\Delta_6 = -19.0$
6	38	34	31	10	11	0		$D(7, A) = 42$	$D(7, B) = 22.2$	$\Delta_7 = -19.8$
7	42	36	36	13	17	9	0			

$A = \{1, 3\}$	$D(7, A) = 42$	$D(7, B) = 22.2$	$\Delta_7 = -19.8$
$B = \{2, 4, 5, 6, 7\}$			

	1	2	3	4	5	6	7		
1	0							$D(2, A) = 8.5$	$D(2, B) = 29.5$
2	10	0						$D(4, A) = 25.5$	$D(4, B) = 13.2$
3	7	7	0					$D(5, A) = 25.5$	$D(5, B) = 15.0$
4	30	23	21	0				$D(6, A) = 34.5$	$D(6, B) = 16.0$
5	29	25	22	7	0			$D(7, A) = 39.0$	$D(7, B) = 18.75$
6	38	34	31	10	11	0			
7	42	36	36	13	17	9	0		

$A = \{1, 3\}$

$B = \{2, 4, 5, 6, 7\}$

	1	2	3	4	5	6	7			
1	0							$D(2, A) = 8.5$	$D(2, B) = 29.5$	$\Delta_2 = 21.0$
2	10	0						$D(4, A) = 25.5$	$D(4, B) = 13.2$	$\Delta_4 = -12.3$
3	7	7	0					$D(5, A) = 25.5$	$D(5, B) = 15.0$	$\Delta_5 = -10.5$
4	30	23	21	0				$D(6, A) = 34.5$	$D(6, B) = 16.0$	$\Delta_6 = -18.5$
5	29	25	22	7	0			$D(7, A) = 39.0$	$D(7, B) = 18.75$	$\Delta_7 = -20.25$
6	38	34	31	10	11	0				
7	42	36	36	13	17	9	0			

$A = \{1, 3, 2\}$

$B = \{ \text{X} 4, 5, 6, 7 \}$

Updated Cluster
 $A = \{1, 3, 2\}$
 and
 $B = \{4, 5, 6, 7\}$

	1	2	3	4	5	6	7
1	0						
2	10	0					
3	7	7	0				
4	30	23	21	0			
5	29	25	22	7	0		
6	38	34	31	10	11	0	
7	42	36	36	13	17	9	0

$$D(4, A) = 24.7 \quad D(4, B) = 10.0 \quad \Delta_4 = -14.7$$

$$D(5, A) = 25.3 \quad D(5, B) = 11.7 \quad \Delta_5 = -13.6$$

$$D(6, A) = 34.3 \quad D(6, B) = 10.0 \quad \Delta_6 = -24.3$$

$$D(7, A) = 38.0 \quad D(7, B) = 13.0 \quad \Delta_7 = -25.0$$

$$A = \{1, 3, 2\}$$

$$B = \{4, 5, 6, 7\}$$

All differences are negative. The process would continue on each subgroup separately.

Recent Hierarchical Clustering Methods

- ▶ Major weakness of agglomerative clustering methods
 - ▶ do not scale well: time complexity of at least $O(n^2)$, where n is the number of total objects
 - ▶ can never undo what was done previously
- ▶ Integration of hierarchical with distance-based clustering
 - ▶ BIRCH (1996): uses CF-tree and incrementally adjusts the quality of sub-clusters
 - ▶ ROCK (1999): clustering categorical data by neighbor and link analysis
 - ▶ CHAMELEON (1999): hierarchical clustering using dynamic modeling

Cluster Analysis

1. What is Cluster Analysis?
2. A Categorization of Major Clustering Methods
3. Partitioning Methods
4. Hierarchical Methods
5. **Density-Based Methods**
6. Summary



Density-Based Clustering Methods

- ▶ Clustering based on density (local cluster criterion), such as density-connected points
- ▶ Major features:
 - ▶ Discover clusters of arbitrary shape
 - ▶ Handle noise
 - ▶ One scan
 - ▶ Need density parameters as termination condition
- ▶ Worst case complexity of $O(m * m)$, where m is number the objects. Space requirement is $O(m)$

Density-Based Clustering: Basic Concepts

- ▶ Two parameters:

- ▶ **Eps**: Maximum radius of the neighbourhood

- ▶ **MinPts**: Minimum number of points in an Eps-neighbourhood of that point

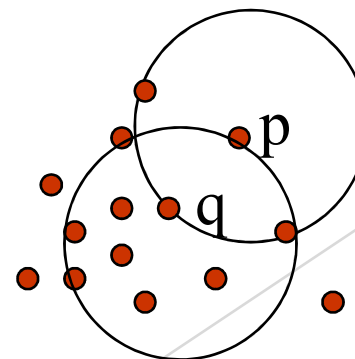
- ▶ $N_{Eps}(p): \{q \text{ belongs to } D \mid \text{dist}(p, q) \leq Eps\}$

- ▶ **Directly density-reachable**: A point p is directly density-reachable from a point q w.r.t. Eps , $MinPts$ if

- ▶ p belongs to $N_{Eps}(q)$

- ▶ core point condition:

$$|N_{Eps}(q)| \geq MinPts$$



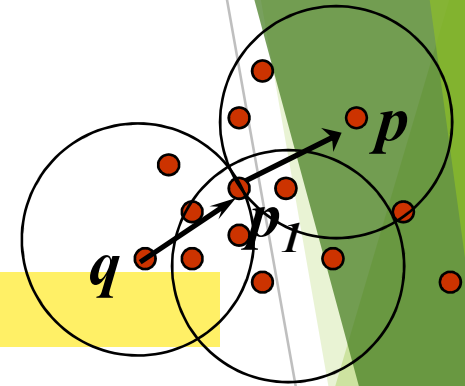
MinPts = 5

Eps = 1 cm

Density-Reachable and Density-Connected

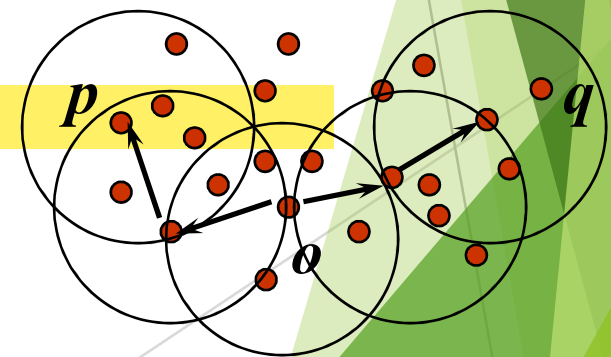
► Density-reachable:

- A point p is **density-reachable** from a point q w.r.t. Eps , $MinPts$ if there is a chain of points p_0, \dots, p_n , $p_0 = q$, $p_n = p$ such that p_{i+1} is directly density-reachable from p_i



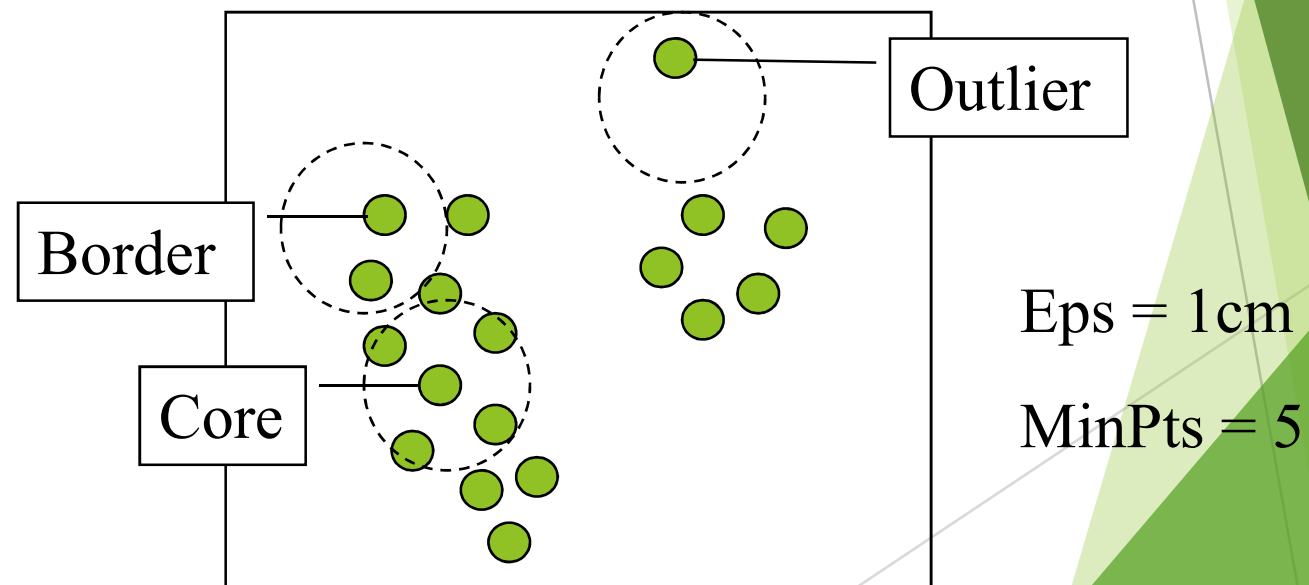
► Density-connected

- A point p is **density-connected** to a point q w.r.t. Eps , $MinPts$ if there is a point o such that both, p and q are density-reachable from o w.r.t. Eps and $MinPts$



DBSCAN: Density Based Spatial Clustering of Applications with Noise

- ▶ Relies on a *density-based* notion of cluster: A *cluster* is defined as a maximal set of density-connected points
- ▶ Discovers clusters of arbitrary shape in spatial databases with noise



DBSCAN: The Algorithm

- ▶ Arbitrary select a point p
- ▶ Retrieve all points density-reachable from p w.r.t. Eps and $MinPts$.
- ▶ If p is a core point, a cluster is formed.
- ▶ If p is a border point, no points are density-reachable from p and DBSCAN visits the next point of the database.
- ▶ Continue the process until all of the points have been processed.

```

DBSCAN(D, eps, MinPts)
  C = 0
  for each unvisited point P in dataset D
    mark P as visited
    NeighborPts = regionQuery(P, eps)
    if sizeof(NeighborPts) < MinPts
      mark P as NOISE
    else
      C = next cluster
      expandCluster(P, NeighborPts, C, eps, MinPts)

expandCluster(P, NeighborPts, C, eps, MinPts)
  add P to cluster C
  for each point P' in NeighborPts
    if P' is not visited
      mark P' as visited
      NeighborPts' = regionQuery(P', eps)
      if sizeof(NeighborPts') >= MinPts
        NeighborPts = NeighborPts joined with NeighborPts'
  if P' is not yet member of any cluster
    add P' to cluster C

regionQuery(P, eps)
  return all points within P's eps-neighborhood (including P)

```

If Epsilon is 2 and minpoint is 2, what are the clusters that DBScan would discover with the following 8 examples: $A1=(2,10)$, $A2=(2,5)$, $A3=(8,4)$, $A4=(5,8)$, $A5=(7,5)$, $A6=(6,4)$, $A7=(1,2)$, $A8=(4,9)$.

Apply DBSCAN algorithm and illustrate the discovered clusters.

What if Epsilon is increased to 10 ?

- Find the Epsilon neighborhood of each point?

$N2(A1)=\{\}$; $N2(A2)=\{\}$; $N2(A3)=\{A5, A6\}$; $N2(A4)=\{A8\}$;

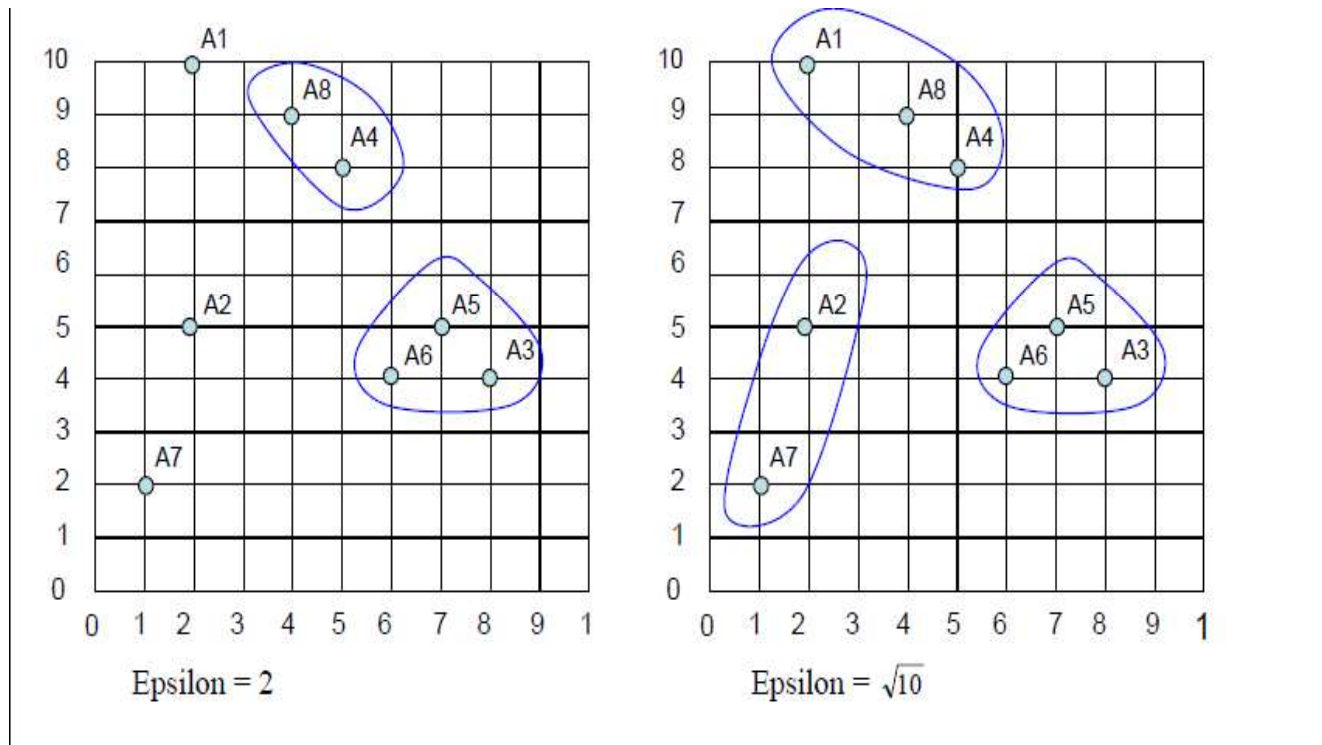
$N2(A5)=\{A3, A6\}$;

$N2(A6)=\{A3, A5\}$; $N2(A7)=\{\}$; $N2(A8)=\{A4\}$

So $A1$, $A2$, and $A7$ are outliers, while we have two clusters $C1=\{A4, A8\}$ and $C2=\{A3, A5, A6\}$

- If Epsilon is 10 then the neighborhood of some points will increase:
- $A1$ would join the cluster $C1$ and $A2$ would joint with $A7$ to form cluster $C3=\{A2, A7\}$.

Clusters : DBSCAN



DBSCAN: Sensitive to Parameters

Figure 8. DBScan results for DS1 with MinPts at 4 and Eps at (a) 0.5 and (b) 0.4.

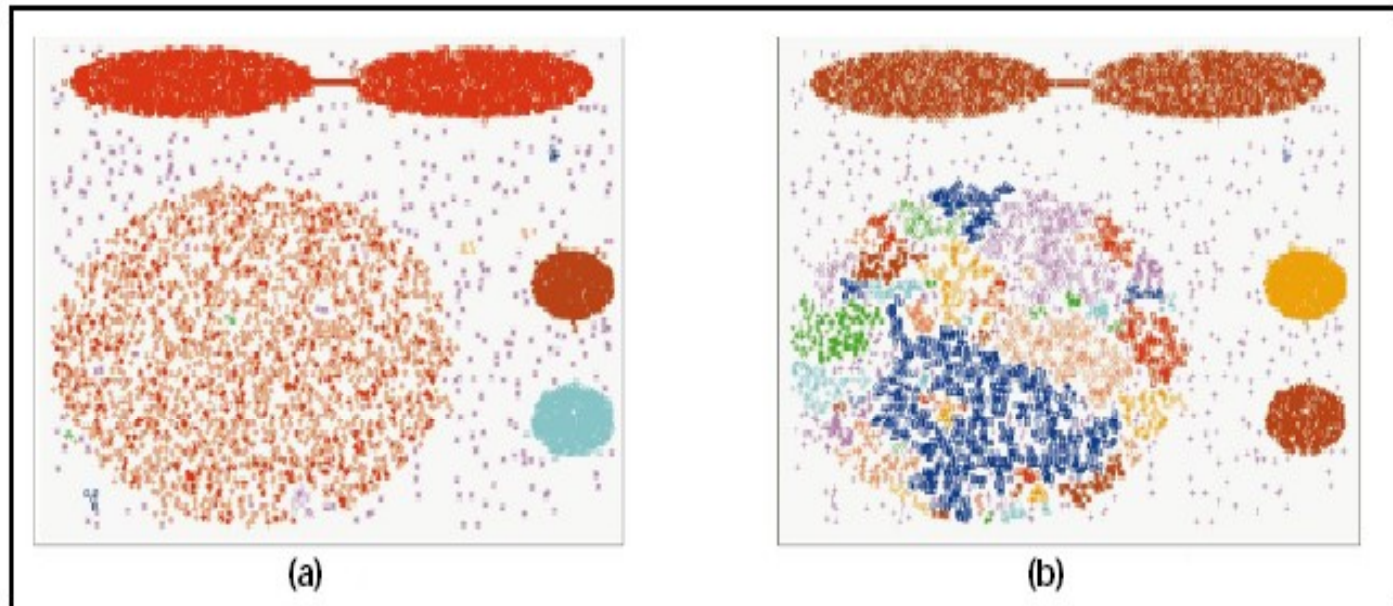
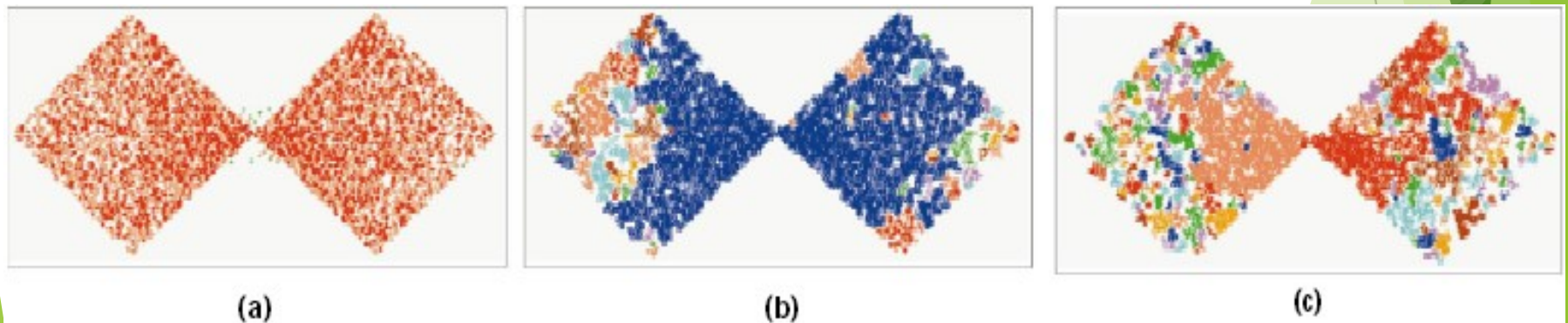


Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



Clustering High-Dimensional Data

- ▶ Clustering high-dimensional data
 - ▶ Many applications: text documents, DNA micro-array data
 - ▶ Major challenges:
 - ▶ Many irrelevant dimensions may mask clusters
 - ▶ Distance measure becomes meaningless—due to equi-distance
 - ▶ Clusters may exist only in some subspaces
- ▶ Methods
 - ▶ Feature transformation: only effective if most dimensions are relevant
 - ▶ PCA & SVD useful only when features are highly correlated/redundant
 - ▶ Feature selection: wrapper or filter approaches
 - ▶ useful to find a subspace where the data have nice clusters
 - ▶ Subspace-clustering: find clusters in all the possible subspaces
 - ▶ CLIQUE, ProClus, and frequent pattern-based clustering

References

- [1] Tan, Pang-Ning, Michael Steinbach, and Vipin Kumar. *Introduction to data mining*. Pearson Education India, 2016.
- [2] Han, Jiawei, Jian Pei, and Micheline Kamber. *Data mining: concepts and techniques*. Elsevier, 2011.
- [3] Lecture Notes by Raymond Wong <http://www.cse.ust.hk/~raywong/comp5331/>
- [4] Batra, Aishwarya. "Analysis and Approach: K-Means and K-Medoids Data Mining Algorithms." *ICACCT, 5th IEEE International Conference on Advanced Computing & Communication Technologies*. 2011.
- [5] Dunham, Margaret H. *Data mining: Introductory and advanced topics*. Pearson Education India, 2006.
- [6] http://cis.csuohio.edu/~sschung/CIS660/chapter_2_SecondJHan.pdf